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.

# Modeling RNA Evolution: In-Silico and In-Vivo

Part 1: RNA Evolution in the RNA World.

Part 2: The Evolution of a RNA Virus: RS Virus in New Zealand

A thesis presented for the degree of

Master of Science

in

BioMathematics

at Massey University, Palmerston North New Zealand

James William Matheson 2004

Copyright © 2004 James William Matheson

#### Abstract

We look at two aspects of the evolution of RNA.

First we look at RNA replication dynamics in an early RNA world context. Experimental evidence (Spiegelman *et al.* 1965, Biebricher *et al.* 1981) shows that under some conditions RNA evolves towards small quickly replicating molecules. We investigate what conditions are sufficient for a population of RNA molecules to evolve towards a balanced population of molecules. This is a population not completely dominated by a single length of molecule. We consider two models: A linear model in which indel rate is inversely proportional to length and a game theory model in which reproductive efficiency depends on the distribution of molecule lengths within a population (this is linked to catalytic efficiency). Models are investigated using analytic, numerical and simulation methods. The linear model is not sufficient to support a population with balanced length distribution. Simulation methods show that the game theory model may support such a population.

We next look at RNA evolution in the context of RNA virus evolution. Using virus samples taken over a thirty year period we investigate the evolution of Respiratory Syncytial Virus (RSV) in New Zealand. RSV most strongly affects infants and the elderly, causing cold like symptoms in mild cases and bronchiolitis or occasionally death in severe cases. New Zealand has a higher incidence of RSV bronchiolitis per head of population than many other developed countries. We compare New Zealand strains of the virus to those isolated overseas to investigate if New Zealand may have significantly different strains. We look at the evolution of the virus within New Zealand looking for evidence of antigenic drift, as well as analysing substitution rates and selection at individual codon sites. No evidence is found to suggest that New Zealand has significantly different strains of RSV from other countries. We conclude the higher rate of severe RSV in New Zealand must be caused by factors other than virus strain. The portion of the virus analysed shows strong evidence of being under positive selective pressure. This and other analyses suggest that RSV may be undergoing antigenic drift.

### Acknowledgements

I would like to say a special thank you to my supervisors, Mike Hendy, David Penny and Barbara Holland for their support, and for the many drafts of this document they have read and corrected. All mistakes that remain are my own.

I would like to thank the Allan Wilson Centre and IFS for their financial support. Thank you to my friends and colleagues there: Klaus, Michael, Tim, Bhalchandra and Barbara for keeping me sane. Thank you to Tim for putting up with so many computer questions and last but not least to Susan and Joy for helping with all the day to day things.

Thank you to Joanna Kirman, Fenella Rich and Catherine Cohet at the Malaghan Institute, without whom I would not have been able to do the RSV section of this thesis. Thank you for putting up with my ignorance of immunology and for sequencing all those snot samples.

Thank you also to Dr. Keith Grimwood of Wellington Hospital and Sue Huang of ESR for your part in the RSV project.

#### Thanks go to:

- Wim Hordijk for help with the initial simulations.
- Igor Boglaev and Andreas Dress for your help with the differential equations description of the game theory.
- Allen Rodrigo for taking the time to tell me about your work.
- Brett Ryland for suggesting a bilinear interpolation.

• James Chai, for putting up with my many enquiries about Helix.

Finally thanks go to my family for being so supportive, especially in my transition from PhD to Masters.

James W. Matheson 5th of December 2004.

### Preface

This project initially started as a PhD thesis and was transmuted into a Masters thesis after one year for personal reasons. In its initial stages the project focused on RNA replication dynamics in an RNA world situation. Simulations for this study were done on the Helix supercomputing cluster.

I also wanted to do some more practical work so when data became available at the Malaghan Institute on Respiratory Syncytial Virus (RSV) I agreed to analyse it as part of a project investigating RSV in New Zealand lead by Dr. Joanna Kirman of the Malaghan Institute.

The majority of the analyses in both parts of this project were done using R (R Development Core Team 2004). High speed simulation code was written in ANSI C. Parallel code used the MPI Parallel library.

**Motivation:** RNA, with its ability to encode both genotype (sequence) and phenotype (folding) in the same molecule is thought to have preceded DNA and protein as a carrier of genetic information by some scientists. An RNA world, in which RNA is the primary information carrying and catalytic molecule, is postulated to have been the first stage of evolution. The plausibility and structure of such a world rests on how RNA behaves. This motivates the study of early RNA evolution.

The study of RSV is motivated primarily by its medical significance. RSV has it largest effect in infants where it causes cold like symptoms in mild cases and bronchiolitis and occasionally death in severe cases. RSV data over a long time period (30 years) is not often available, and there is little data available on New Zealand RSV, so the discovery of New Zealand RSV samples from 1967 to 1997 in an ESR freezer offers a useful opportunity to study this virus.

## Contents

Abstractii
Acknowledgementsiii
Preface
Table of Contents
Figure List
Table Listx
Glossary of Terms and Abbreviationsxi

1	Intr	oducti	on	1
	1.1	What	is RNA?	2
2	The	Evolu	tion of Biological Molecules.	4
	2.1	Indel I	Bate Models	5
		2.1.1	Base Model Specification	7
		2.1.2	The Ur-Si model.	11
		2.1.3	Nonuniform Replication: The Nr-Si model	14
		2.1.4	The Ur-Ai model	15
		2.1.5	The Nr-Ai model	19
		2.1.6	Conclusion: Indel Rate Models	20
	2.2	Game	Theory	21
		2.2.1	Evolutionary Game theory.	21
		2.2.2	The Evolutionary Game Model: The Simple Model	24
		2.2.3	Expanding the game: The Full Model	27
		2.2.4	Simulating the Full Model	29
		2.2.5	The Planar Games.	30

		2.2.6	The Bilinear Game	35
		2.2.7	Proving the hypothesized ESS	41
		2.2.8	Differential Equation Model of the Bilinear Game	42
		2.2.9	Conclusion: Game Theory Models	44
	2.3	Conclu	usion: All tested models	44
3	Res	pirato	ry Syncytial Virus (RSV)	45
	3.1	The R	S Virus: An Introduction	46
	3.2	Gener	al Methods	47
		3.2.1	Networks	52
		3.2.2	Permutation Tests	60
	3.3	Evolut	tion of RSV in New Zealand	63
		3.3.1	Introduction	63
		3.3.2	Methods	64
		3.3.3	Results	77
		3.3.4	New Zealand Conclusions.	93
	3.4	Intern	ational comparisons.	94
		3.4.1	Introduction	94
		3.4.2	Methods	94
		3.4.3	Results	97
		3.4.4	International Conclusions.	106
	3.5	Final	Conclusions	106
4	$\operatorname{Ref}$	erence	s	108
5	App	pendix	1: Sequence Group and Number Keys.	111
	5.1	Numb	er Key for figures 3.14 and 3.15.	111
	5.2	Group	bings for Figure 3.16: RSV B F-Protein Alignment.	111
	5.3	Group	bings for Figure 3.17: RSV B G-Protein Alignment.	112
	5.4	Numb	per Key for Figures 3.18 and 3.19	112
	5.5	Group	bings for Figures 3.21 and 3.22: RSV A G-Protein Alignment.	114

## List of Figures

2.1	Steady State distribution when long molecules have a smaller chance of	
	indels in the Ur-Si model.	13
2.2	Steady state of the the Nr-Si model with nonuniform replication compared	
	to the Ur-Si model	15
2.3	The effect of $\varepsilon$ on the steady state	18
2.4	Steady state of a highly asymmetric system under the Nr-Ai model	20
2.5	Example of Multiple ESS in a 2 player 3 strategy game	24
2.6	Intersection Types.	31
2.7	Graphs showing the surfaces generated by the two subtypes of the game	
	matrix $G$	32
2.8	Simulation results after $10^4$ cycles using the concave game matrix starting	
	from a uniform distribution.	33
2.9	Simulation results after $10^4$ cycles using the convex game matrix starting	
	from a uniform distribution.	34
2.10	Four steps of a bilinear interpolation.	36
2.11	Surface generated by the bilinear game matrix	38
2.12	Simulation results using the bilinear game matrix with random seed $s_1$ after	
	1 million cycles	39
2.13	Simulation results using the bilinear game matrix with random seed $s_2$ after	
	1 million cycles	40
2.14	Steady state distribution of molecules generated by equation (2.15)	43
3.1	Cartoon of the RS Virus	46
3.2	RSV B Sampling dates	48
3.3	RSV A sampling dates	48
3.4	Section of RSV Genome under study	49
3.5	sUPGMA tree for RSV B G718 primer site data.	51
3.6	Splits on a tree	54
3.7	Incompatible Splits	55

3.8	Example Network.	56
3.9	Treeness Scale.	57
3.10	All Splits in the RSV B G718 Alignment	58
3.11	Medium Number of Splits Displayed.	59
3.12	Least Number of Splits Displayed	60
3.13	Antigenic Drift Effects Example	64
3.14	RSV B in New Zealand G718 Data	78
3.15	Local clustering summary results for RSV B in New Zealand	80
3.16	F-Protein Amino Acid 1 to 55	82
3.17	RSV B G-Protein Amino Acid 259 to 299 (N-terminus)	83
3.18	RSV A In New Zealand FG-A Network.	85
3.19	Local clustering to summary results for RSV A in New Zealand.	87
3.20	Graph of Number of substitutions vs Number of years separating samples	
	for RSV A	88
0.01		
3.21	New Zealand RSV A G-Protein amino acids 156 to 225	90
3.21 3.22	New Zealand RSV A G-Protein amino acids 156 to 225	90 91
3.21 3.22 3.23	New Zealand RSV A G-Protein amino acids 156 to 225	90 91 95
<ol> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> </ol>	New Zealand RSV A G-Protein amino acids 156 to 225	90 91 95 96
<ol> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> <li>3.25</li> </ol>	New Zealand RSV A G-Protein amino acids 156 to 225	90 91 95 96 98
<ol> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> <li>3.25</li> <li>3.26</li> </ol>	New Zealand RSV A G-Protein amino acids 156 to 225	90 91 95 96 98 00
<ul> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> <li>3.25</li> <li>3.26</li> <li>3.27</li> </ul>	New Zealand RSV A G-Protein amino acids 156 to 225	90 91 95 96 98 00
<ul> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> <li>3.25</li> <li>3.26</li> <li>3.27</li> <li>3.28</li> </ul>	New Zealand RSV A G-Protein amino acids 156 to 225	<ul> <li>90</li> <li>91</li> <li>95</li> <li>96</li> <li>98</li> <li>00</li> <li>01</li> <li>02</li> </ul>
<ul> <li>3.21</li> <li>3.22</li> <li>3.23</li> <li>3.24</li> <li>3.25</li> <li>3.26</li> <li>3.27</li> <li>3.28</li> <li>3.29</li> </ul>	New Zealand RSV A G-Protein amino acids 156 to 225	<ul> <li>90</li> <li>91</li> <li>95</li> <li>96</li> <li>98</li> <li>00</li> <li>01</li> <li>02</li> <li>04</li> </ul>

## List of Tables

2.1	Model Type Definitions.	6
2.2	The four types of indel rate model investigated	7
2.3	A summary of the important features of the indel rate models, $\ldots$ , $\ldots$	8
2.4	Transition probabilities for a molecule of length $l$ under the Ur-Si model. $\ .$	12
2.5	Transition probabilities for a molecule of length $l$ under the Ur-Si model. $\ .$	16
2.6	Transition probabilities for a molecule of length $l$ under the Ur-Ai model	
	with parameter $\varepsilon$ ,	16
2.7	The Payoff matrix for prisoners P1 and P2 from the point of view of P1 in	
	the prisoners dilemma game	22
2.8	Game Payoffs for the template molecule in the simple game	25
2.9	Numerical values of payoff to template molecule in the simple game model.	27
2.10	Names of payoffs for the simple game and the reaction type they correspond	
	to	28

#### LIST OF TABLES

## Glossary of terms and abbreviations.

Asymmetric Indel rate
The rate at which insertions occur does not equal the rate at which deletions occur.
Asymmetry Parameter
Determines the relative probability of insertions and deletions.
Bilinear Interpolation
Estimation of internal points of a surface by linear interpolation from two boundaries.
Cluster Measure
Measures the degree of homogeneity of a sequence.
Copy Error
The rate at which copy errors occur during replication given as a probability per replica-
tion of an error occurring.
Cytodomain
The part of a protein inside the viral capsid.
Degree of Congruence
Extent to which split partitions coincide with partitions defined by identical property
values.
Degree of Homogeneity
Measure of the degree to which terms of a sequence are identical.
<b>DNA</b>

Deoxyribonucleic Acid. Molecule that encodes genetic information of all cellular life forms.

Ectodomain
The part of a protein outside the viral capsid (exposed to environment).
<b>ESS</b>
Evolutionary Stable Strategy: An ESS is a strategy such that if a population is using it
then that population is not vulnerable to invasion by mutants using any other strategy.
Evolutionary Game theory
Modelling the evolutionary process using game theory.
<b>FG-B</b>
Nested PCR Primer site on RSV B. Covers portions of G-Protein ectodomain.
<b>FG-A</b>
Nested PCR Primer site on RSV A. Covers portions of the G-Protein ectodomain.
<b>F-protein</b>
Fusion Protein (Surface protein that facilitates cell entry).
G714
RT-PCR Primer site on RSV B. Covers portions of G-Protein ectodomain, intergenic
spacer region and F-Protein cytodomain.
<b>C718</b> 40
BT-PCB Primer site on RSV B. Covers portions of G-Protein ectodomain intergenic
spacer region and F-Protein cytodomain
opaosi rogion and i i rotom cytodomani.
G-protein
Attachment Glycoprotein (Surface protein that facilitates cell attachment).

Incompatible Splits
Two or more splits that can not be displayed on the same tree.
Indel
The addition or deletion of bases from a molecule during replication.
Indel rate
Rate at which Indels occur during replication. Length dependent.
Mixed strategy
Individuals within a population use strategy 1 some of the time and strategy 2 the re-
mainder of the time.
Nonuniform Replication
Replication rate is dependent on a molecules length
Replication rate is dependent on a molecules length.
Nr-Ai Model
Nonuniform-Replication Asymmetric-Indel rate Model.
Nr-Si Model
Nonuniform-Replication Symmetric-Indel rate Model.
Path Length
Number of edges separating two vertices connected by a path in a connected graph.
Delumente de structeres
Polymorphic strategy
Individuals within a population consistently use one strategy. Some individuals use strat-
egy 1 while the remaining individuals use strategy 2.

#### LIST OF TABLES

Property Function
Function returning sequence of property values for the terms of the input sequence.
Replication rate
Used to determine the number of offspring produced by a molecule. Length dependent.
<b>RNA</b>

Ribonucleic Acid. A single stranded biomolecule.

Split	. 52
A bipartition of a set or sequence .	

Split Congruence
When a Splits partitions coincide with partitions defined by terms of the sequence having
dentical property values.

sUPGMA	50
Serial UPGMA	

Symmetric Indel rat	e			6
For a molecule of given	length, th	ie rate insertions	occur equals the	rate deletions occur.

### LIST OF TABLES

Uniform Replication
Replication rate is uniform over all molecules irrespective of their length.

Ur-Ai Model	. 7
Uniform-Replication Asymmetric-Indel rate Model.	
Ur-Si Model	. 7
Uniform-Replication Symmetric-Indel rate Model.	

### 1 Introduction

This thesis investigates several aspects of the evolution of biological molecules. RNA is the molecule we concentrate on. We take a brief look at RNA in section 1.1.

In Section 2 we look at how RNA might have been involved in the early development of life. Eigen (1971) considers constant length self replicating biomolecules (specifically RNA). Eigen's model goes from early chemical self organization in evolution to the formation of hypercycles, which are self-replicating autocatalytic cycles. There is some argument over the veracity of his claims for hypercycles as the path to more complex biomolecules (Boerlijst and Hogeweg 1991, Zintzaras *et al.* 2002), however our area of interest lies before the formation of hypercycles in the area of RNA evolution. We look at variable length self replicating biomolecules and ask the question: Can we find a set of conditions sufficient to create a stable population of molecules with a balanced length distribution? This question arises from experimental observations (Spiegelman *et al.* 1965, Biebricher *et al.* 1981) showing that, under some conditions, RNA will evolve towards a highly biased (short) length distribution. This unlike what is seen in nature today. We attempt to answer this question by using mathematical models and simulation methods to investigate RNA replication dynamics in different model systems.

In section 3 we move from the theoretical to the practical implications of RNA evolution. We look at the evolution of Respiratory Syncytial Virus (RSV) in New Zealand. RSV is common in infants. Mild infection causes symptoms similar to a common cold, severe infection can cause bronchiolitis and death. New Zealand has a higher incidence of hospital admissions from RSV bronchiolitis than many other developed countries (Vogel *et al.* 2003). We investigate if this is due to New Zealand having different strains of RSV to other countries, as well as looking at the characteristics of the virus's evolution in New Zealand.

#### 1 INTRODUCTION

#### 1.1 What is RNA?

RNA stands for **R**ibo**N**ucleic **A**cid which is a single stranded cousin to DNA. RNA like DNA is constructed from a linear chain of nucleotides attached to a sugar phosphate backbone. Unlike DNA, RNA is usually single stranded and the nucleotides used are Adenine, Cytosine, Guanine and Uracil (not Thymine) often abbreviated to A, C, G and U.

Due to their single stranded nature RNA molecules are free to fold back on themselves. This means that nucleotides that fold to be adjacent can form hydrogen bonds (base-pair) creating 2D structures such as loops and hairpins. Different nucleotides form different strengths of bond. The strongest bond is  $C \equiv G$  followed by A = U. There is also weak binding affinity in the bond C - A. Bases with strong base pairing affinities are said to complement each other, C and G are complementary base pairs as are A and U. The 2D structures formed by RNA can in turn fold in 3D space to form complex structures. Some of these 3D structures will provide the chemical binding sites which allow RNA to catalyse chemical reactions.

Modern theories of the origin of life assume an RNA-world stage (Yarus 1999). This is a stage of evolution that is dominated by RNA. In these theories founding populations of RNA molecules are produced by natural RNA synthesis from nucleotides on ancient earth. RNA replication is aided by catalytic RNA called ribozymes. The founding population gradually evolves, by mechanisms such as that discussed in Eigen (1971), towards the production of protein and eventually DNA. The situation in the RNA-world differs from the modern situation; in the RNA-world RNA was both catalyst and information carrying molecule. In the modern situation the functions of catalyst and information carrier are separated between protein and DNA respectively. RNA catalysts for RNA processing (ribozymes) are essential to the RNA-world hypothesis as they form the basis for theories of self-replicating systems of RNA molecules. Though the existence of efficient ribozymes is yet to be experimentally proven, there is reason to believe it will be, with groups such as that of David Bartel (Lawrence and Bartel 2003) finding molecules that have good RNA polymerase activity but limited processivity (their ability to catalyse other molecules is

#### 1 INTRODUCTION

not long lived).

In a hypothetical RNA world environment containing free nucleotides and ribozymes some RNA sequences can undergo replication (Spiegelman *et al.* 1965). This process involves a complementary copy of the molecule being created from the original by pairing each base in the original unwound strand with its complement. This complement is complemented in turn to create a replica of the original sequence. There are no known error correcting mechanisms in this process so RNA replication is prone to errors. Errors can take the form of miscoded bases (called 'substitutions') or the addition or deletion of bases from the molecule (these are called insertions and deletions respectively and are collectively referred to as indels).

In this thesis we look at a formal model of aspects of the early RNA world (section 2) as well as how RNA, in the form of RNA viruses, evolves in the world today (section 3).