

# Computational Detection and Analysis of Polyadenylation Signals

Ashwin Hajarnavis

Darwin College  
&  
The Wellcome Trust Sanger Institute

March 2005

A dissertation submitted for the degree of  
Doctor of Philosophy  
at the University of Cambridge

**Preface:**

This dissertation is the result of my own work and includes nothing, which is the outcome of work done in collaboration except where specifically indicated in the text.

No part of this thesis is being submitted for any other qualification or at any other University.

**Acknowledgements:**

I should like to thank my supervisor, Richard Durbin, for his constant advice and guidance. The Sanger Institute has been an incredible place to work, on account of the many people I have met there. Numerous though they are, some deserve a special mention. Firstly, Kevin Howe and Ian Korf, for sharing so much of their time, and discussing everything from *C. elegans*. Also Lachlan Coin and Thomas Down, for being inspirational sources of technical help. Access to data was made simple, thanks to help from all at WormBase, in particular Daniel Lawson and Keith Bradnam. I am grateful to Sam Griffiths-Jones and Alex Bateman for many ideas, and for the trip into the world of RNA. Finally, current and former members of our lab; Diego DiBernardo, Marc Sohrmann, Irmtraud Meyer, David Carter, Avril Coghlan, and Mark J Minichiello.

This work was funded by the Medical Research Council and The Wellcome Trust.

## Summary:

Many computational techniques exist for the prediction of genes from genome sequence, and for their functional characterisation. Less well understood, however, are the sequences that cause processing and regulation of these genes. One such sequence is the polyadenylation signal, which is required for the expression of most eukaryotic genes. The ability to detect polyadenylation signals accurately means that genomes can be annotated to a greater extent. Although this can be carried out in the laboratory, a computational method is much faster and cheaper, especially considering the acceleration in the sequencing of whole genomes.

A particular gain is that a polyadenylation signal prediction also provides a predicted end to the untranslated region (UTR) lying downstream (3') of a gene's stop codon. This region can contain regulatory motifs, which can dictate properties such as when, where, and how a gene is expressed. Knowledge of gene regulation is as important as gene function if we are to try and gain a full understanding of systems biology from genome sequencing.

In this thesis, I present the development of a piece of software for detecting sequence signals in genome sequence.

I then develop a model for the polyadenylation signal in the nematode worm *Caenorhabditis elegans* and show that the predictions are accurate, leading to the publication of good quality 3' UTR data sets.

Models are then built for three other species, and a comparison made with existing methods.

A comparison between polyadenylation signals of *C. elegans* and the closely related *C. briggsae* follows, which leads onto the discovery of a putative regulatory motif, conserved between the ribosomal protein 3' UTR sequences of both species.

# Contents

1.	An Introduction to 3' Ends and Polyadenylation Signals.....	1
1.1.	Preamble.....	1
1.2.	Overview of untranslated region molecular biology .....	4
2.	PAjHMMA – Parameter Adjustable Java Hidden Markov Model Architecture .	18
2.1.	Introduction .....	18
2.2.	An overview of hidden Markov models (HMMs).....	18
2.3.	Software design .....	25
2.4.	Conclusion.....	39
3.	A Probabilistic Model for 3' End Formation in <i>C. elegans</i> .....	41
3.1.	Introduction .....	41
3.2.	Background.....	42
3.3.	Model building .....	43
3.4.	Model evaluation .....	57
4.	Polyadenylation Signal Prediction in Other Eukaryotes .....	70
4.1.	Introduction .....	70
4.2.	Data Acquisition .....	72
4.3.	Nucleotide Frequencies.....	73
4.4.	Model testing.....	84
4.5.	Discussion .....	89
4.6.	Conclusions .....	92
5.	On the Evolution of 3'UTRs and Polyadenylation Signals .....	94
5.1.	Introduction .....	94
5.2.	Conservation of absolute position .....	95
5.3.	Polyadenylation signals in aligned orthologues .....	97
5.4.	Discussion – On the evolution of polyadenylation signals .....	107
6.	Concerning a Sequence Element Detected in Ribosomal mRNAs.....	109
6.1.	Introduction .....	109
6.2.	Background.....	109
6.3.	Model building .....	113
6.4.	Model testing.....	120
6.5.	Results.....	120
6.6.	Discussion .....	124
6.7.	Conclusions .....	125
6.8.	Collaboration – the analysis of another 3' UTR binding motif.....	126
7.	Conclusions.....	128
	References .....	132
	Appendices .....	139