# Computational detection of non-coding RNAs in genomes

Yen-Hua Huang

This dissertation is submitted for the degree of Doctor of Philosophy

The Wellcome Trust Sanger Institute and Churchill College, Cambridge

# Declaration

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.

The work in this thesis has not been submitted in whole, or in part, for a degree, diploma, or any other qualification at any other university.

Yen-Hua Huang

March 2008, Cambridge, UK

# Abstract

Noncoding RNAs (ncRNAs) have become implicated in a variety of regulatory mechanisms as well as structural roles, suggesting that functional ncRNAs may be more prevalent in genomes than previously supposed. Nonetheless, *in silico* ncRNA finding is difficult, even though a mass of genome sequence is publicly available. Few computational approaches are really reliable for genome-wide ncRNA finding. This thesis is devoted to assessing available approaches and trying new solutions for finding ncRNAs in genomes.

In the first half of this thesis, reasons that may contribute to the slow progress of genome-wide ncRNA finding are explored. A comprehensive analysis on a genome-wide scale of the credibility of currently used signals for classifying ncRNAs is conducted. Two factors, conservation of ncRNAs in human-mouse syntenic regions and abundance of covariations between human-mouse synteny-conserved ncRNAs, are evaluated. The result reveals that current comparative-genomics-based methods may not be able to find ncRNAs effectively in mammalian genomes. In addition, possible genomic features that could distinguish real ncRNAs from pseudogenes are investigated. Two different criteria, distribution of bit scores and physical clustering in genomes, are applied to filter out tRNA pseudogenes and to enrich *bona-fide* tRNA genes. Physiological roles of the tRNA genes in human-mouse synteny-conserved clusters are discussed and the degradation patterns of tRNA pseudogenes are analyzed.

In the second half of this thesis, computational techniques are applied to model signals that may be potentially useful for genome-wide ncRNA finding. A sparse Bayesian learning algorithm, Eponine, is applied to model the transcription start sites of mammalian ncRNA genes that are transcribed by RNA polymerase III. In addition to modelling *cis*-regulatory elements for transcription, a new computational module, which extends the capability of

Eponine to learn motifs consisting of both primary sequences and RNA secondary structures, is created. The capability of this new module is demonstrated by applying it to analyze several known cases of ncRNA motifs. The strength and the weakness of applying this new computational approach for finding ncRNAs are discussed.

# Acknowledgements

# Table of Contents

# List of Tables

viii

# List of Figures