# Computational localization of promoters and transcription start sites in mammalian genomes

Thomas Down

This dissertation is submitted for the degree of Doctor of Philosophy

Wellcome Trust Sanger Institute and Sidney Sussex 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.

Thomas A. Down

May 2003, Cambridge, UK

# Abstract

A number of large genomes have now been sequenced, and biologists are now faced with the challenge of identifying all the functional pieces of sequence, and understanding how they contribute to the development and life of the organism. While identification of protein coding genes, and annotating their products, has been progressing well, there are a great many open questions relating to the regulatory regions which control the expression of these genes.

Here, I investigate the question of identifying and annotating promoters, one of the most important regulatory signals in the genome, which mark the points where transcription is initiated, and regulate the transcription of genes. I present a new computational method, EponineTSS, which can predict transcription start sites in bulk genomic sequence data with excellent sensitivity and specificity. Unlike the existing methods, it gives an indication of the actual location of the transcription start site. Comparisons with available experimental data suggest that the positional accuracy of these predictions is very good. Results from this method are included as part of the Ensembl human genome annotation.

Having located transcription start sites for genes, I also discuss the use of results from comparative genomics the estimate the extent of the functional promoter region upstream of the start site. I show that the extent of promoters is very variable, and that promoter size is correlated with the function of the gene for whose regulation it is responsible. Genes associated with developmental processes tend to have particularly large, and thus presumably complex, promoters, with the homeobox transcription factors among the most extreme examples.

I also introduce sparse Bayesian learning, a recently developed approach to supervised machine learning which can be applied to the training of a wide range of model types, and embodies the principle of selecting the simplest possible model to explain the observed data. I demonstrate a new technique which makes sparse Bayesian learning much more scalable,

allowing it to be applied to very large and complex problems, and present a convenient, freely available Java library which provides a general-purpose implementation of this technique. This library was used here in the training of the transcription start site predictor, but has a wide range of applications in computational biology and beyond.

# Acknowledgments

Thanks to the Wellcome Trust, and everyone at the Sanger institute who helped and encouraged me as I explored this fascinating area of research. I particularly appreciate the efforts of my project supervisor, Dr. Tim Hubbard, who has helped me along the path towards this thesis, accepted the many digressions, and added a valuable voice of caution and rigour when analyzing new results. Just because you're paranoid, it doesn't mean the data *isn't* out to get you.

All the programs I wrote during this project rest on the strong foundation of the BioJava library. I have to thank everyone who has ever contributed to this project: documentation, code, bug reports and ideas are all vital to the success of an open source effort. Right from the outset, Matthew Pocock has been a guiding light for the project. He's also a great coder, target for trying out new ideas, and friend.

Throughout this all, I've enjoyed the love, support, and good humour of Emily White, who has kept me going through the difficult moments and helped celebrate the joyous ones. Thank you.

# Contents

# List of Figures

x

# List of Tables