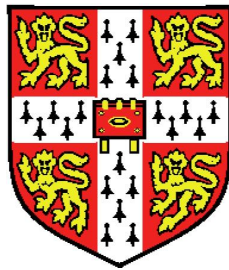# Inference and classification of eukaryotic *cis*-regulatory motifs

Matias Piipari

Wellcome Trust Sanger Institute

University of Cambridge

A thesis submitted for the degree of

*Doctor of Philosophy*

6th September, 2010

With ♡ to Kaisa.

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.

# Acknowledgements

I would like to thank my supervisors Dr Tim Hubbard and Dr Thomas Down for the valuable advice and support I have been given during my project. I am also grateful to members of my PhD committee: Dr Derek Stemple, Dr Alex Bateman and Dr Jurg Bahler. I am deeply honoured to have studied my four years in Cambridge with the generous Wellcome Trust studentship.

My sincere thanks go to the Hubbard research group members Markus Brosch, Mutlu Dogruel, Jenny Mattison, and especially Daniel James who read the preliminary versions of this dissertation. My time at Sanger has been made both happy and productive thanks to many of the other Sanger Institute future doctors. I would especially like to thank my dear friends Alexandra Nica, Leopold Parts, Sergei Manakov, Steve Pettitt and Marija Buljan for the many conversations, lessons and laughs they have let me experience.

I will always be grateful for the love and endless encouragement for intellectual pursuit that I have received from my parents, brother and sister. Finally, I am indebted to the love and friendship of my dear wife Kaisa.

Matias Piipari, September 2010, Cambridge, UK.

# Abstract

Regulation of gene expression by networks of sequence specific transcription factors is one of the most important control mechanisms that defines the expression pattern of a genome. Describing transcriptional regulatory networks requires a near complete knowledge of the transcription factors present in the cell, as well the DNA binding sites to which each of the TFs is able to bind. Recent years have witnessed advances in both directions. High coverage transcription factor annotations have become available for many sequenced eukaryotic genomes. Improvements have also been made in profiling DNA specificity motifs for eukaryotic transcription factors, *in vitro* and *in vivo*.

The theme of my work has been the application and development of computational methods for inferring regulatory motifs from promoter sequence, and finding clues to the function of computationally inferred DNA motifs. Functional annotation of inferred motifs led me to conduct a comparative study of the familial relationships between regulatory motifs, the conclusion of which was a probabilistic motif family model I call the 'metamotif'. The metamotif, I will show, allows improved prediction of the DNA binding domain family for *de novo* inferred motifs, and is an effective way of encoding prior information about known DNA binding domain families to a motif inference algorithm. The use of familial prior information improves the sensitivity to detect regulatory motifs contained in the large promoter sequences that are common to higher eukaryotic genomes. The metamotif guides motif inference towards types of sequence signal that are expected *a priori* to be present in the sequence set of interest, thereby improving and supplementing traditional regulatory motif inference algorithms.

I have also assessed several published *de novo* DNA motif inference algorithms by challenging them to infer a complete set of regulatory motifs from a large series of *Saccharomyces cerevisiae* promoters. This work provides a novel way to assess performance of regulatory motif inference methods, and is made possible by the availability of an experimentally determined regulatory motif dictionary for the *S. cerevisiae* genome. In addition to benchmarking motif inference methods compared to a reference motif set, I make use of many of the rich genomics resources available for study of the budding yeast. These include curated lists of TF target genes based on ChIP-chip and gene expression studies of wild type and knockout yeasts, a close-to-complete list of TF motif from the JASPAR database, and a 7-way sequence conservation score across the genome, as well as sequence variation data from the *Saccharomyces* Genome Resequencing Project.

Development of sensitive regulatory motif inference algorithms continues to be important in gaining understanding of eukaryotic gene regulation by sequence specific transcription factors. In particular I believe that methods that integrate different sources of biological evidence, such as metamotifs, gene expression and ChIP-seq, to sequence motif inference will be highly important to the field.

# Contents

# List of Figures