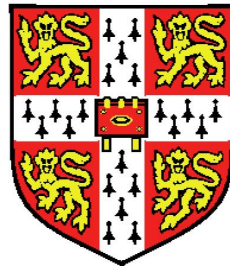# Motif based computational identification of protein subcellular localisation

Mutlu Doğruel

Wellcome Trust Sanger Institute

University of Cambridge

A thesis submitted for the degree of

*Doctor of Philosophy*

January 2008

To my deeply missed father,

Sami Han Doğruel.

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 express my gratitude to Dr. Tim Hubbard for being a great supervisor. His support, invaluable suggestions and comments during our weekly meetings are much appreciated. I am also deeply indebted to my PhD committee members Dr. Ewan Birney and Dr. Ian Dunham for their time and advice. I would also like to thank my examiners Prof. Søren Brunak and Dr. Nick Goldman for reading the thesis thoroughly and for their useful feedback on the preliminary version of the thesis.

My sincere thanks go to Dr. Thomas Down for his useful suggestions and also patience during the discussions we made (usually accompanied with a cup of coffee) on Biojava, NestedMICA, Eponine and many other topics... Many thanks also to the other members of Tim's research group at the Sanger Institute: Dr. Andreas Prlić, and will-be Drs. Markus Brosch, Matias Piipari, and Jenny Mattison for making my PhD days more enjoyable, and being excellent friends. I think Markus did a great job by encouraging me to join the Cambridge University gliding club. Being in the air over the weekends in the summer was a good way to relax and get ready for the productive work days;)

Finally, I am deeply and forever indebted to my parents, sisters and brother for their love, support and encouragement throughout my entire life: Sizi çok seviyorum!

Mutlu Doḡruel, January 2008, Cambridge, UK.

# Abstract

Discovering overrepresented patterns in amino acid sequences is an important step in protein functional annotation which includes the identification of subcellular localisation. I adapted and extended NestedMICA, an *ab initio* protein motif finder originally developed for finding transcription binding site motifs, to find short protein signals, and compared its performance with another popular protein motif finder, MEME.

In order to assess NestedMICA as a protein motif finder, I have tested it on synthetic datasets produced by spiking instances of known motifs from protein databases into a randomly selected set of protein sequences. Apart from the artificially implanted motifs, NestedMICA also successfully recovered subcellular localisation signals from biologically-authentic test sets. NestedMICA found most of the short test protein motifs spiked into a test set of sequences at different frequencies. In all the assessment experiments, its overall motif discovery performance was better than that of MEME.

As a practical application of NestedMICA, I developed a novel Support Vector Machines based protein subcellular classification tool,

Lokum, for eukaryotic protein subcellular localisation prediction, covering all major localisation classes for animal, fungal and plant sequences. It uses targeting and retention signal motifs reported by NestedMICA, and other protein features including transmembrane topologies and amino acid composition. Additionally, in Lokum I use bipartite nuclear localisation signals obtained by adding protein support to Eponine, a tool originally developed for transcription start site modeling. Lokum does not use sequence similarity, or any other *a priori* knowledge such as known nuclear localisation signals by searching databases.

I compared proteins targeted into the nuclei and nucleoli in terms of the features used in Lokum, and also their predicted disorder regions. I demonstrate that it is possible to computationally distinguish these two sub-nuclear protein categories.

Finally, as an alternative to the transmembrane topology predictor TMHMM that is used in Lokum, I designed and tested a new prototype program that is based on hidden Markov models (HMM). The HMM has been trained by a novel, nested sampling based transition probability optimisation procedure.

# Contents

# List of Figures

# List of Tables