# Fine-Mapping and Functional Analyses of Genetic Variants Driving Local Adaptations in Humans

University of Cambridge

Corpus Christi College



A thesis submitted for the degree of

*Doctor of Philosophy*

Michał Szpak

The Wellcome Trust Sanger Institute
Wellcome Genome Campus
Hinxton, Cambridge
CB10 1SA, UK

September 2016

*To my parents and siblings,*
*my granny and auntie AM*

# Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work carried out under the supervision of Prof. Chris Tyler-Smith at the Wellcome Trust Sanger Institute, while member of Corpus Christi College, University of Cambridge, and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 60,000 words excluding bibliography, figures, tables, equations and appendices.

Michał Szpak
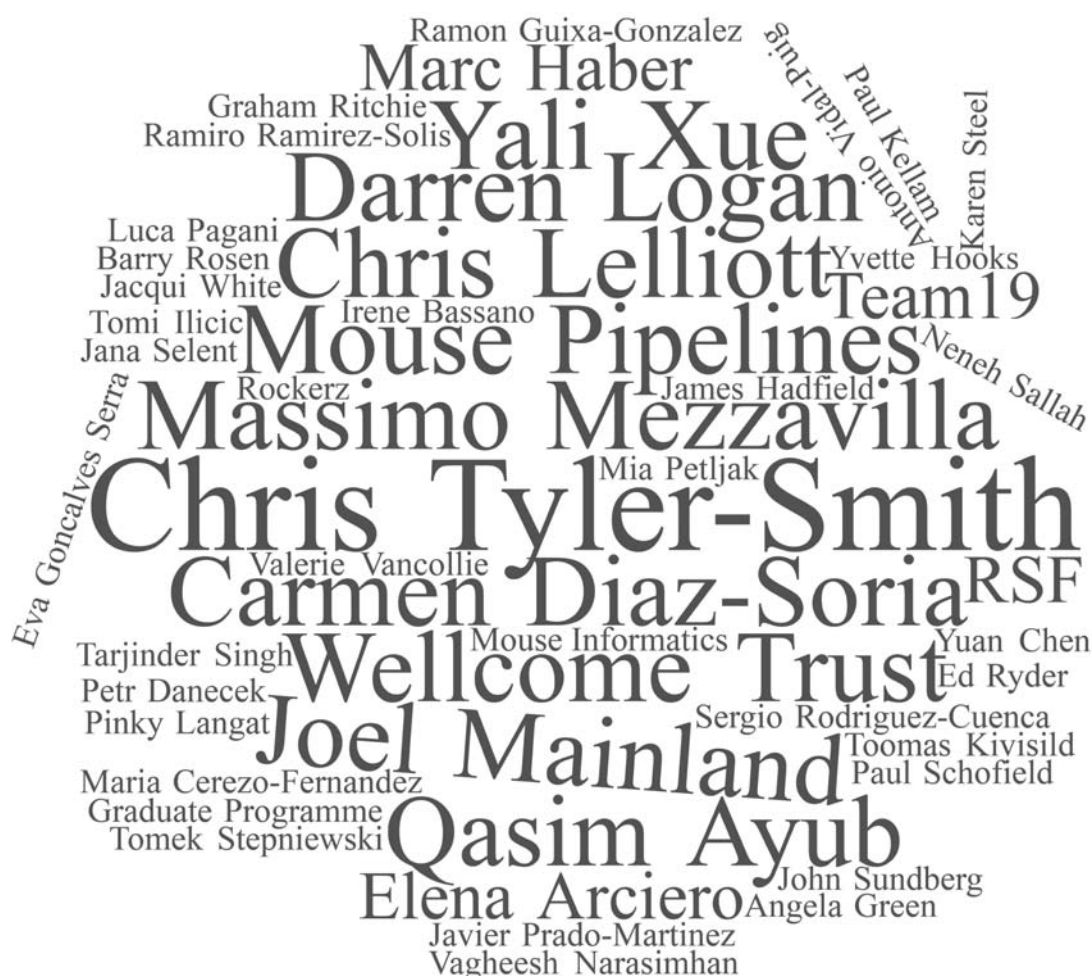
September 2016

7

# Acknowledgements



Figure 0. Acknowledgements. First and foremost, I would like to thank my supervisor Prof. Chris Tyler-Smith for giving me the opportunity to carry out this project and for all his invaluable advice, support and encouragement. Many thanks also to my secondary supervisors Dr Yali Xue and Dr Qasim Ayub for their critical and constructive assessment of my work. I would also like to thank current and former team19 members for their continual advice and guidance, especially Dr Massimo Mezzavilla, Dr Marc Haber, Dr Javier Prado-Martinez, Elena Arciero, Yuan Chen, Vagheesh Narasimhan, Dr Luca Pagani and Dr Maria Cerezo-Fernandez, who contributed to this project. I also thank my internal and external collaborators including Carmen Diaz-Soria, Dr Irene Bassano, Prof. Paul Kellam, Dr Joel Mainland, Dr Darren Logan, Dr Toomas Kivisild, Dr Petr Danecek, Dr Graham Ritchie and Lukasz Szpak for providing their expertise and fruitful collaborations. I also thank the WTSI Research Support Facility, Mouse Informatics and Mouse Pipelines, and in particular Dr Chris Lelliott, Dr Ramiro Ramirez-Solis, Dr Jacqui White, Dr Barry Rosen, Dr Ed Ryder, Yvette Hooks, Valerie Vancollie, Dr Angela Green and Hannah Wardle-Jones for productive discussions, generation and phenotyping of the mouse models and free dissemination of the data used in this project. I am also very grateful to my thesis committee members Dr Paul Schofield and Dr Darren Logan, WTSI Graduate Programme and the Committee of Graduate Studies, as well as the Wellcome Trust for my PhD studentship. I also wish to acknowledge external collaborators who joined our efforts in future development of this project, namely Prof. Karen Steel, Dr John Sundberg, Dr Sergio Rodriguez-Cuenca, Prof. Antonio Vidal-Puig, Tomek Stepniewski, Dr Ramon Guixa-Gonzalez and Dr Jana Selent. Lastly, a special thanks to my fellow PhD students, especially Tomi, Carmen, Neneh, Mia, Pinky, TJ, James and Eva, as well as the members of the 'Murray's breakfast table' and my non-Sanger friends for their support and being the best.

# Abstract

The genetic basis of human evolutionary adaptation and the resulting population diversification has been of great interest. A common approach has been to scan genomes for population-genetic signatures of positive selection, yielding vast lists of thousands of candidates. Here, we first took advantage of these data to perform a meta-analysis of published selection screens and assessed their concordance using a Selection Support Index (*SSI*) which weights, combines and evaluates signals of selection on a per-gene basis. Our analysis revealed both the low overall agreement of previous genome-wide selection scans and some strong candidates. The focus of positive selection studies in humans thus needs to move from candidate locus discovery to pinpointing underlying causal variants and further investigation of their biological significance. We developed a new computational method for this, Fine-Mapping of Adaptive Variation (*FineMAV*), which combines population differentiation, derived allele frequency and a measure of molecular functionality to prioritise candidate selected variants for functional follow-up. We calibrated and tested *FineMAV* using eight 'gold standard' examples of experimentally-validated causal variants underlying positive selection, and were able to pick out the known functional allele in all instances. We used this approach to identify the best candidate variants driving local adaptations in the 1000 Genomes Project Phase 3 SNP dataset including Africans, admixed Americans, Europeans, and East and South Asians. *FineMAV* top hits were overall enriched for high *SSI* scores, and we also report many novel examples, including rs6048066 in *TGM3* associated with curly hair and rs7547313 in *SPTA1* associated with erythrocyte shape and possibly malaria resistance in Africa, as well as rs201075024 in *PRSS53* linked to hair shape in South Asia. We extended our analyses to additional populations including Egyptians, Ethiopians, Greeks, Lebanese and non-admixed Native Americans, picking up interesting hits in Peruvian Quechua and Ethiopian Gumuz in genes involved in immunity and energy metabolism. The highest scoring *FineMAV* variant in Native Americans was rs34890031 in *LRGUK* associated with spermatogenesis. We then performed functional follow-up on chosen candidates. Our *in vitro* studies focused on comparison of the ancestral and derived forms of the

*OR10H3* olfactory receptor, and of *FUT2* involved in susceptibility to viruses, but were limited by technical issues. We also investigated the functions of six genes showing strong signals of selection using mouse knock-outs. The curly vibrissae (whiskers) of *Prss53* knock-out mice supports our hypothesis of selection in *PRSS53* due to hair shape in humans, while *Herc1* knock-out mice show a range of abnormalities affecting hearing, blood plasma chemistry and energy metabolism. Finally, we initiated the generation of nine mouse knock-ins carrying a human selected allele, which will be subjected to future collaborative phenotyping, focusing on hair shape, reproduction, energy metabolism and hearing as appropriate. Our work is thus facilitating the identification of causative alleles driving human adaptations.

# Table of contents

14

# List of figures

18

# List of tables

# List of equations