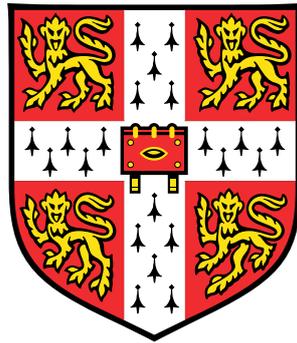# Identification of germline variants that predispose to familial melanoma

**Aravind Sankar**

Wellcome Sanger Institute

University of Cambridge

This dissertation is submitted for the degree of

*Doctor of Philosophy*

St Edmund's College

2020

To Appa, Amma and Aditya: without whom this would not be possible.

# 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 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. This dissertation contains less than 60,000 words and has less than 150 figures, exclusive of tables, footnotes, bibliography, and appendices.

Aravind Sankar
2020

# Acknowledgements

There is an old adage which goes "It takes a village to raise a child". I believe it also takes a village to obtain a doctoral degree. I would like to take this opportunity to thank the members of my village who helped me get to where I am today.

First and foremost, I would like to thank my supervisor Dave Adams for his endless support, guidance, wisdom and kindness. I wouldn't know where to begin thanking you for all the things you have done for me so I would just like to say - thank you for everything.

I would also like to thank my other supervisor, Vivek Iyer, who mentored, trained and tutored me with the patience of a saint. It is due to him that I can confidently call myself a bioinformatician today.

I am extremely grateful to Professor Tim Bishop from the University of Leeds for taking the time to e-mail, phone and make trips to Cambridge to teach me statistics.

I would like to acknowledge the Wellcome Sanger Institute and the MELGEN network for funding my PhD and for providing me with a platform to collaborate, interact and learn from current and future world-leaders in research.

I would like to thank every member of the GenoMEL consortium that provided samples, ideas, inputs and suggestions over the last 4 years. This project would not have been possible without them.

To Sofia - Thank you for starting and finishing your PhD with me, for the dinners, the conference trips, for always lending an ear and for being my voice of reason.

To Gemma - Thank you for being my friend, my driver, my flat mate and my confidante. Every day got a little bit easier with you around.

To Nicky - Thank you for the tea trips, the joint birthday celebrations and for reminding me to always look on the bright side of life.

To Marco - Thank you for taking the effort to make me feel at home in Cambridge on my first day and for continuing to do that ever since.

To Daniela - Thank you for always being there even when you are in Mexico.

To Katharina - Thank you for your positivity, for encouraging me to run and for being a social butterfly.

To Annie - Thank you for taking the time to remind me that I cannot take my time with my thesis!

To everyone else in team113 - Thank you for making me feel excited to come to work everyday. I will see you at lunch at 12.

To Prakaash and Mahathi - Thank you for sticking with me.

Finally, to my family. Appa, Amma, Paati and Aditya - For encouraging me to fly while being my safety net, for being living examples of showing that hard work pays off and for your unconditional love and support - Thank you.

# Abstract

Melanoma is an extremely aggressive malignancy with a poor prognosis in advanced disease. While GWAS and exome analysis have helped to identify loci linked to the development of the disease, these studies have explained predisposition to melanoma in only a fraction of cases. Thus, the majority of the genetic factors that contribute to the pathogenesis of melanoma are yet to be defined. This project aims at identifying novel genes and pathways involved in the development of familial melanoma, and also identify loci which predispose individuals to disease development.

308 individuals from 133 different families previously diagnosed with melanoma were sequenced through a mixture of exome or whole genome sequencing. Multiple workflows were established to analyse the dataset for novel driver mutations. A novel approach of combining association and linkage analysis was established for the variants in the coding region to identify genes with high burden of mutations where the variants segregated with the disease within the pedigrees. The role of non-coding variants and structural variants in melanoma onset was also investigated through additional workflows in the whole-genome sequenced individuals.

Non-synonymous mutations were found in *CDKN2A, BRCA1, POT1* and *BAP1*. Disruptive variants were also observed in novel genes such as *EXO5, TP53AIP* and *AMER1*. An increased burden on variants in transcription factor binding motifs were observed in genes including *SYK* and *SRC*. A large deletion upstream of *CDKN2A* was identified. Genes including *ATR* and *FAT1* were identified to have a higher burden of disruptive variants that segregated with the disease within the cases through the novel combined association-linkage analysis.

Disruptive germline variants that could play a role in familial melanoma development were identified in multiple genes through a combination of several approaches.

# Contents

# List of Figures

# List of Tables

# Nomenclature

**Roman Symbols**

CGC   Cancer Gene Census

ExAC   Exome Aggregation Consortium

GC   Count of individuals for each genotype

GC NFE   Count of Non-Finnish European individuals for each genotype

gnomAD   Genome Aggregation Database

GQ   Genotype Quality

GRCh37   Genome Reference Consortium Human Build 37

GRCh38   Genome Reference Consortium Human Build 38

gVCF   Genomic Variant Calling Format

HGNC   HUGO Gene Nomenclature Committee

OR   Odds Ratio

PRS   Polygenic Risk Scores

TF   Transcription Factors

TFBM   Transcription Factor Binding Motifs

VCF   Variant Calling Format

VEP   Variant Effect Predictor