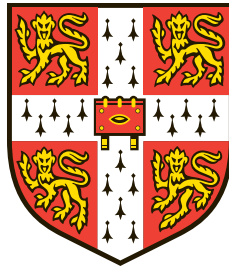# Host and pathogen genetics associated with pneumococcal meningitis

## John Andrew Lees

Wellcome Trust Sanger Institute
Jesus College, University of Cambridge

July 2017

This dissertation is submitted for the degree of
Doctor of Philosophy

# Declaration

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in the Preface and specified in the text.

It is not substantially the same as any that I have submitted, or, is being concurrently submitted for a degree or diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. I further state that no substantial part of my dissertation has already been submitted, or, is being concurrently submitted for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text.

It does not exceed 60 000 words in length, as required by the School of Biological Sciences.

John Andrew Lees

July 2017

# Acknowledgements

My main thanks have to go to Stephen Bentley, who has supervised me in exactly the way I would have wished. As well as being a steady hand on the tiller of my PhD, he has also shown me that it is possible to have a well-adjusted life, be kind to others and still be a great researcher. I will never forget all that he taught me (about Special Brew). Similar thanks must also go to my other supervisors. Jeff Barrett, who had no obligation to do so, fully accepted me into his research group. His group meetings and my conversations with him have really shaped this research and the way I think. Julian Parkhill has always been helpful and available for answering those tricky bacterial genomics questions no-one else could, and has given me many opportunities to present my work to other researchers (often in sunnier climes). Thanks too to Carl Anderson and John Welch, the other members of my thesis committee, whose freely contributed ideas have made this work better than it otherwise would have been.

Many collaborators have made this work possible. Nick Croucher is a pneumococcal master who I am lucky to know. Jukka Corander I am likewise lucky to know, and while explicitly involved with SEER I have felt his influence throughout this PhD. My Dutch friends Diederik, Matthijs, Philip, Arie and Bart have worked very hard on the unique dataset at the core of this thesis, and have been a pleasure to collaborate with. A special mention to Philip who even saw fit to invite me to his wedding. Paul Turner helped me throughout with the Maela dataset, and organised an excellent meeting in Cambodia without which chapter three wouldn't exist.

To the people who I have talked to about my PhD, thanks for your contributions (tangible and intangible) to everything I've done here: Sumana, Katie, Tom, Jeremy, Liam, Marcia, Theresa, Simon, James, Sophia, Becca, Leo, Izzy, Claire, Darryl and Alison, and all the other members of teams 81 and 143 past and present. Thanks too to all of pathogen informatics and the graduate office, who have doubtless helped me in many ways. Finally, I am very grateful to the Wellcome Trust and MRC for funding this research.

In the knowledge that this is the only page most readers of this document will look at, the pressure to be witty or memorable is greatest here. I guess you'll have to live with the Special Brew reference.

# Summary

## Host and pathogen genetics associated with pneumococcal meningitis

John Andrew Lees

Meningitis is an infection of the meninges, a layer of tissue surrounding the brain. In cases of pneumococcal meningitis (where the bacterium *Streptococcus pneumoniae* is the causative agent) this causes severe inflammation, requiring intensive care and rapid antibiotic treatment. The contribution of variation in host and pathogen genetics to pneumococcal meningitis is unknown. In this thesis I develop and apply statistical genetics techniques to identify genomic variation associated with the various stages of pneumococcal meningitis, including colonisation, invasion and severity.

I start by describing the development of a method to perform genome-wide association studies (GWAS) in bacteria, which can find variation in bacterial genomes associated with bacterial traits such as antibiotic resistance and virulence. I then applied this method to longitudinal samples from asymptomatic carriage, and found lineages and specific variants associated with altered duration of carriage. To assess meningitis versus carriage samples I applied similar analysis techniques, and found that the bacterial genome is crucial in determining invasive potential. As well as bacterial serotype, which I found to be the main effect, I discovered many independent sequence variants associated with disease. Separately, I analysed within host-diversity during the invasive phase of disease and found it to be of less relevance to disease progression.

Finally, I analysed host genotype data from four independent studies using GWAS and heritability estimates to determine the contribution of human sequence variation to pneumococcal meningitis. Host sequence accounted for some variation in susceptibility to and severity of meningitis. The work concludes with a combined analysis of pairs of bacterial and human sequences from meningitis cases, and finds variation correlated between the two.

# Contents

# List of Figures

9

# List of Tables

# Acronyms

**AF**  allele frequency. 56, 57

**AIC**  Akaike information criterion. 84, 85

**ALF**  artificial life framework. 58, 71

**AMP**  anti-microbial peptide. 21, 26

**BAM**  binary sequence alignment/map. 111, 112

**BAPS**  Bayesian analysis of population structure. 48, 54, 59, 61–63, 79, 180, 193

**BFGS**  Broyden–Fletcher–Goldfarb–Shanno. 65, 66

**CDS**  coding sequences. 138, 139, 142, 145

**CFU**  colony forming unit. 135

**CI**  confidence interval. 48, 118, 126

**CMH**  Cochran–Mantel–Haenszel. 48, 54, 75, 79, 193

**CNV**  copy number variant. 39, 108, 112, 124, 143

**COG**  cluster of orthologous genes. 30, 46, 49, 55–57, 112, 120, 122, 188, 193

**CPP**  closest phylogenetic-pairs. 118

**CSF**  cerebrospinal fluid. 17–21, 23, 77, 108–110, 114, 117, 132, 134–152, 185, 186, 192, 196

**CSV**  comma separated values. 176

**d.f.**  degrees of freedom. 36, 56, 66

**DSM**  distributed string mining. 55, 56, 71, 73

**FWER**  family-wise error rate. 36, 67

**GoNL** The Genome of the Netherlands. 162

**GOS** Glasgow outcome score. 20, 118

**GTR** generalised time reversible. 58, 60

**GWAS** genome wide association study. 17, 33, 36–39, 41–52, 54, 55, 57, 63, 75, 77, 79, 81, 82, 98, 106, 108–112, 116, 119, 121, 125, 135, 147, 151, 152, 154–156, 161, 167–172, 183, 185–187, 189, 191, 193–197

*H. influenzae* *Haemophilus influenzae*. 20, 24, 155

**HLA** human leukocyte antigen. 43, 154, 175, 182

**HMM** hidden Markov model. 82, 84, 85, 95, 189

**HPD** highest posterior density. 147, 148

**HRC** haplotype reference consortium. 162, 163

**HWE** Hardy-Weinberg equilibrium. 158, 160, 162

**ICE** integrative conjugative element. 29, 31, 74, 91, 126, 129, 130

**ICU** intensive care unit. 156

**IPD** invasive pneumococcal disease. 18, 24

*ivr* inverting variable restriction. 32, 116, 117, 119, 131, 136, 146, 147, 252

**JC** Jukes-Cantor. 60

**KC** Kendall-Colijn. 60–62

*L. monocytogenes* *Listeria monocytogenes*. 20, 47, 58, 155

**LD** linkage disequilibrium. 30, 34–38, 42, 44–46, 49, 57, 73–75, 79, 88, 96, 98–100, 102, 161, 162, 164, 166, 177, 196

**LMM** linear mixed model. 39, 50, 86, 88–90, 93, 99, 102, 105, 120, 164, 187–189, 191, 193, 194, 246

**LOD** logarithm of odds. 33

**LoF** loss of function. 28, 39, 51, 124, 125, 128–131, 140, 141, 151, 191

**LRT** likelihood ratio test. 62, 66, 67, 90, 118, 164, 187

***M. tuberculosis*** *Mycobacterium tuberculosis*. 43, 46, 47, 50, 128, 195

**MAC** membrane attack complex. 22, 27, 112

**MAF** minor allele frequency. 34, 36, 38, 39, 42, 55, 68, 71, 77, 96, 98, 99, 124, 128, 156, 158, 160–163, 165, 166, 170, 175–178, 180, 183

**MCMC** Markov-chain Monte Carlo. 116, 118, 132, 163

**MDS** multidimensional scaling. 63–65, 67, 68, 119, 176

**MIC** minimum inhibitory concentration. 93

**MLST** multi-locus sequence typing. 30, 47, 59, 61, 62, 108, 139, 143

**MNP** multiple nucleotide polymorphism. 110

**MRCA** most recent common ancestor. 58, 194

***N. gonorrhoeae*** *Neisseria gonorrhoeae*. 66

***N. meningitidis*** *Neisseria meningitidis*. 20, 21, 43, 46, 47, 99, 109, 135, 136, 138, 139, 142–146, 149, 150, 152, 155

**NCD** normalised compression distance. 60–62

**NJ** neighbour joining. 60–62

**NT** non-typable. 25, 31, 82, 85, 86, 90, 95

**OR** odds-ratio. 19, 45, 48, 49, 71, 72, 128, 165, 166, 170, 175, 178, 180, 183

**OU** Ornstein-Uhlenbeck. 118

***pbp*** penicillin binding protein. 29, 49

**PCA** principal component analysis. 39, 115, 158, 178, 251

**PCR** polymerase chain reaction. 132, 147

**PCV** pneumococcal conjugate vaccine. 21, 31, 82, 195

**PEER** probabilistic estimation of expression residuals. 178–180

***ply*** pneumolysin. 26, 195

**QC** quality control. 36, 109, 111, 155, 157, 160, 162, 163, 176