# Functional and evolutionary analyses of pneumococcal genome variation
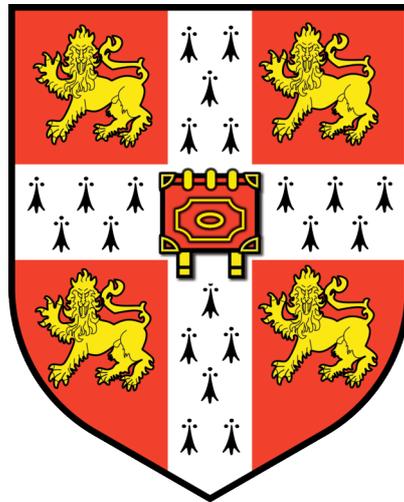
**Nicholas Jason Croucher**

**Downing College**

**University of Cambridge**

**2011**

**This dissertation is submitted for the degree of Doctor of Philosophy**

**Declaration**

I hereby declare that this dissertation is my own work and contains nothing that is the outcome of work done in collaboration with others, except where specifically indicated here or elsewhere in the text.

All sequence data were produced by the sequencing and research and development teams at the Wellcome Trust Sanger Institute. The annotation of complete genome sequences was performed in conjunction with Dr. Stephen Bentley (Wellcome Trust Sanger Institute, Cambridge). Phylogenetic analyses of bacterial genome sequences were performed in collaboration with Dr. Simon Harris (Wellcome Trust Sanger Institute, Cambridge). All microarray experiments and statistical analyses were performed by the BμG@S group (St. George's Hospital, London). All statistical analyses of Omnilog data were performed by Lars Barquist (Wellcome Trust Sanger Institute, Cambridge).

None of the work described herein has been previously submitted for the purpose of obtaining another degree. This dissertation does not exceed 60,000 words in length, as required by the School of Biological Sciences.

Nicholas Jason Croucher

August 2011

**Acknowledgements**

I thank Julian Parkhill for allowing me the opportunity to undertake this work and for all his help and advice, in particular his ability to patiently correct me when I have been stubbornly wrong. I also owe a huge debt of gratitude to Stephen Bentley, upon whom the burden of my supervision has fallen most heavily; despite this, he has remained unflinchingly positive, understanding and encouraging, and he has proved to be a fantastic, friendly and wise supervisor. Similarly, Nick Thomson has been a perennially enthusiastic and supportive manager; it has been a pleasure to work for all three. I also thank Duncan Maskell for his perceptiveness and advice as my external supervisor, and the other members of my advisory committee: Gordon Dougan, who has also invested much effort in helping with my laboratory work, and Alex Bateman, who has organised the graduate students at the Sanger Institute so effectively.

Much of my informatics work could not have been completed without the help of Simon Harris and Matthew Holden, who have been two of the most friendly, conscientious and helpful people I could possibly hoped to have shared an office with. Thomas Otto also deserves much credit for shouting and singing me into being a better systems user. I am also grateful to Theresa Feltwell and Sally Whitehead for their tolerance of me, and the disruption I inevitably cause, in running their respective laboratories. I thank Maria Fookes and Del Pickard for their vast microbiological expertise and general helpfulness, as well as their unendingly friendly and amusing ways, along with Trevor Lawley, who also taught me a very practical approach to laboratory etiquette. I am also thankful to all of teams 15 and 81 for their patience and understanding, with a special mention of congratulations to the ever entertaining and hospitable Alan 'Aunt' Walker. Of course, none of this work could have been completed without the help of the informatics, systems, sequencing and library making teams of the Sanger Institute, which has been a brilliant place to work.

Externally, I have enjoyed many productive discussions with Bill Hanage and Christophe Fraser of Imperial College, London. Gavin Patterson, of the Vetinary School, Cambridge, was very kind in initially providing me with strains and advice on laboratory techniques. Tim Mitchell, of the University of Glasgow, has been a

friendly and invaluable collaborator. In addition, I am greatly appreciative of the kindness of all our other collaborators who have provided us with relevant strains, which have been essential to this project.

I also thank my family, in particular my mum, my dad and my grandparents who have been so supportive and encouraging throughout these seven years at Cambridge, and beyond. The friends I have made during this time have also been central in making the bad times survivable and the good times so enjoyable. My final thanks are to Alina, who has had to demonstrate more patience and understanding than could reasonably be expected of anyone, and has made me so happy.

p.s. Mum, to save you flicking through to the end: no, I still haven't cured anything.

**Abstract**

**Functional and evolutionary analyses of pneumococcal genome variation**

**Nicholas Jason Croucher**

*Streptococcus pneumoniae* (the pneumococcus) is a human nasopharyngeal commensal and respiratory pathogen responsible for a high burden of morbidity and mortality worldwide. The bacterium's primary virulence factor appears to be its polysaccharide capsule, of which there are more than 90 different serologically-distinguishable types (serotypes). Although this categorisation was originally used for tracing pneumococcal epidemiology, the bacterium is naturally transformable, and hence is able to switch serotypes through horizontal exchange of capsule biosynthesis (*cps*) gene clusters. Therefore, following the emergence of multidrug-resistant lineages in the late 1970s, superior, multilocus-based typing schemes were devised for following pneumococcal evolution. Increasing antibiotic resistance also motivated the development of a heptavalent conjugate polysaccharide vaccine, which targeted seven *S. pneumoniae* serotypes, leading to a decrease in pneumococcal disease. However, this impact has been ameliorated by an increase in disease resulting from replacement by non-vaccine serotypes and switching of *cps* loci by strains previously expressing vaccine serotypes. This thesis describes the application of second-generation sequencing technologies to investigating the mechanisms by which the pneumococcus evolves, especially in response to such clinical interventions.

The first part concerns the Pneumococcal Molecular Epidemiology Network clone 1 (PMEN1) lineage, one of the first multidrug-resistant pneumococcal genotypes to become a worldwide problem. Complete sequencing of the *S. pneumoniae* ATCC 700669 type strain, combined with draft sequencing of a global collection of 240 isolates, quantified the impact of recombination across the chromosome, as well as revealing the diversity of conjugative elements and prophage in the population. The acquisition of antibiotic resistances and the evasion of the conjugate polysaccharide vaccine were both evident in among the strains. *In vitro* transformation experiments, in the same genetic background, were then used to perform a more detailed

investigation of the types of homologous recombination events seen in the global population.

The second part of this dissertation describes the use of RNA sequencing to investigate the functional consequences of genomic variation. A novel method was developed and validated, and, when applied to *S. pneumoniae* ATCC 700669, revealed a family of expressed putative coding sequences that were formed by extended forms of the BOX interspersed repeat. This technique was also applied to two closely related strains of the PMEN31 lineage, both isolated from a single case of disease. This allowed the functional consequences of a small number of distinguishing polymorphisms on the global transcriptome to be ascertained, providing an insight into the level of pneumococcal evolution that can occur within an individual. Sequencing further members of this lineage showed that, although highly successful, this lineage has a much more static genotype than that of PMEN1.

The different mechanisms of pneumococcal genome variation are associated with evolution over different timescales, and in response to different selection pressures, but clearly interact in a number of ways. Hence the use of whole genome sequencing, surveying all the variation throughout the chromosome, will be crucial for greater understanding, and therefore improved control, of this important pathogen.

**Contents**

## List of Figures

## List of Tables

## Abbreviations

| | |
|---|---|
| aa | Amino acid |
| ABC | ATP-binding cassette |
| ATCC | American type culture collection |
| ATP | Adenosine triphosphate |
| BAC | Bacterial artificial chromosome |
| BHI | Brain heart infusion |
| BLAST | Basic local alignment search tool |
| bp | Base pair |
| CBP | Choline-binding protein |
| CC | Clonal complex |
| cDNA | Complementary DNA |
| CDS | Coding sequences |
| CGH | Comparative genome hybridisation |
| CSF | Cerebrospinal fluid |
| CSP | Competence stimulating peptide |
| DNA | Deoxyribonucleic acid |
| dNTP | Deoxynucleoside triphosphate |
| ds | Double stranded |
| DUS | DNA uptake sequence |
| EDTA | Ethylenediaminetetraacetic acid |
| EMBL | European Molecular Biology Laboratories |
| FR | Flanking region |
| GC | Guanine and cytosine |
| GMP | Guanosine monophosphate |
| HIV | Human immunodeficiency virus |
| HMM | Hidden Markov model |
| ICE | Integrative and conjugative element |
| IPD | Invasive pneumococcal disease |
| IS | Insertion sequence |
| LB | Luria broth |
| MCS | Multiple cloning site |

| | |
|---|---|
| MEPS | Minimum efficiently processed segment |
| MGE | Mobile genetic element |
| MITE | Miniature inverted repeat transposable element |
| MLEE | Multilocus enzyme electrophoresis |
| MLST | Multilocus sequence typing |
| MMR | Mismatch repair |
| MR | Mosaic recombination |
| mRNA | Messenger RNA |
| nt | Nucleotide |
| NTP | Nucleoside triphosphate |
| ORF | Open reading frame |
| PBP | Penicillin binding protein |
| PCR | Polymerase chain reaction |
| PCV | Polysaccharide conjugate vaccine |
| PEP | Phosphoenolpyruvate |
| PFGE | Pulsed field gel electrophoresis |
| PMEN | Pneumococcal molecular epidemiology network |
| PPI-1 | Pneumococcal pathogenicity island 1 |
| PTS | Phosphotransferase system |
| PUS | *patAB* upregulatory SNP |
| rDNA | Ribosomal DNA |
| RNA | Ribonucleic acid |
| RPKM | Reads per kilobase per million mapped reads |
| rRNA | Ribosomal RNA |
| RSS | Recombined sequence segment |
| RT-PCR | Reverse transcription PCR |
| RUP | Repeat unit of pneumococcus |
| SGST | Second generation sequencing technology |
| SNP | Single nucleotide polymorphism |
| ss | Single stranded |
| ST | Sequence type |
| TCA | Tricarboxylic acid |
| TIR | Terminal inverted repeats |
| tRNA | Transfer RNA |

TSD ................................................ Target sequence duplication
USS ................................................ Uptake signal sequence
UTR................................................ Untranslated region
UV .................................................. Ultraviolet
VNTR ............................................. Variable number tandem repeat