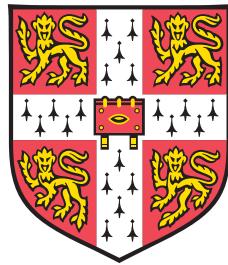# Genomic variation and evolution of *Salmonella enterica* serovars Typhi and Paratyphi A

Kathryn Holt

Wolfson College, University of Cambridge

This dissertation is submitted for the degree of

*Doctor of Philosophy*

August 2009

# Abstract

*Salmonella enterica* serovars Typhi and Paratyphi A are bacterial pathogens that cause typhoid fever in humans. Typhi and Paratyphi A are unusual among *S. enterica* serovars, as they are restricted to systemic infection of humans while most serovars cause gastroenteritis in a broad range of animal hosts. Despite their similarities, Typhi and Paratyphi A are thought to have evolved independently, adapting to the human systemic niche via mechanisms which are still poorly understood. There is little genetic variation within each population, making it difficult to study their evolution or population dynamics.

In this thesis, comparative genomic analysis was used to detect variation within the Typhi and Paratyphi A populations, and to compare the evolution of these two pathogens. A total of 19 complete Typhi genome sequences were compared in order to identify genetic variants, including single nucleotide mutations (SNPs), deletions and insertions of novel DNA. A different approach was taken to study the Paratyphi A population, including the comparison of seven complete genome sequences and development of a novel technique to screen for SNPs in a collection of 160 genomes sequenced in pools. Little evidence was found of selection upon Typhi genes, but there was evidence of diversifying selection in genes coding for the biosynthesis of O-antigen in Paratyphi A. There was evidence in both populations of ongoing accumulation of inactivating mutations which result in loss of gene function. Detailed comparison of this functional gene loss in Typhi and Paratyphi A revealed that many of the same genes were inactivated in both serovars, but the mutations occurred independently and were not the result of horizontal transfer of DNA between their genomes. Comparative analysis of variation in the Typhi and Paratyphi A populations suggested that

Paratyphi A is the younger pathogen, with a most recent common ancestor roughly a third as old as that of Typhi.

Bacteria can harbour plasmids (additional strands of circular DNA) that carry genes encoding resistance to drugs. The plasmids are able to spread between bacterial cells, thereby spreading drug resistance within or between pathogen populations. In this thesis, comparative analysis of plasmid sequences from Typhi and Paratyphi A found that the same type of plasmid was present in both serovars, carrying identical DNA sequences encoding resistance to the drugs used to treat typhoid fever. This demonstrates that the evolution of drug resistance in both serovars is tightly linked. Very closely related sequences were also found in other human bacterial pathogens, highlighting how easily drug resistance can spread.

Single nucleotide variants (SNPs) identified in Typhi and in the drug resistance plasmids were used to develop a high-throughput SNP typing assay with which to study Typhi populations. The SNP typing assay was used to interrogate a global collection of Typhi, as well as local Typhi populations from areas where typhoid is endemic, including regions of Vietnam, Nepal, India and Kenya. The analysis linked strain type with plasmid type for the first time, and demonstrated multiple independent acquisitions of distinct drug resistance plasmids over the past 40 years, culminating in the current dominance of a single plasmid type. Analysis of recent Typhi populations circulating in endemic areas showed that the same Typhi clone now dominates all of these regions, although local diversification has resulted in subtle differences between the populations. Importantly, the dominant Typhi clone was closely associated with the dominant plasmid type, suggesting that the success of the clone and plasmid may have been intimately linked.

# Declaration

This dissertation is my own work and contains nothing
which is the outcome of work done in collaboration with others,
except as specified in the text and Acknowledgements.

The thesis work was conducted from May 2006 to August 2009
at the Wellcome Trust Sanger Institute, Cambridge, UK
under the supervision of
Gordon Dougan (Wellcome Trust Sanger Institute),
Julian Parkhill (Wellcome Trust Sanger Institute), and
Duncan Maskell (Department of Veterinary Medicine,
University of Cambridge).

To my parents,

who introduced me to the world of science,

and to my husband Mike,

who made it possible to stay.

# Acknowledgements

# Contents

# List of Figures

# List of Tables

# Glossary

**bp**     Base pairs

**CDS**     Protein-coding sequences

**contig**     Contiguous sequence assembled from overlapping reads

**Gb**     Gigabase pairs (1 billion bp)

**GTR**     General time reversible substitution model

**homoplasy**     Identity by state but not by descent

**IncHI1**     Plasmid incompatibility type HI1

**indel**     Insertion/deletion mutation

**IS**     Insertion sequence

**IVI**     International Vaccine Instute, Seoul, South Korea

**kbp**     Kilobase pairs (1 thousand bp)

**KEMRI**     Kenya Medical Research Institute, Nairobi, Kenya

**LPS**     Lipopolysaccharide

**Mbp**     Megabase pairs (1 million bp)

**MCMC**     Markov chain Monte Carlo

**MDR**     Multiple drug resistance, defined as resistance to chloramphenicol, ampicillin and co-trimoxazole

**MIC**     Minimum inhibitory concentration, defined as the minimum concentration of an antimicrobial that can inhibit the visible growth of a microorganism

**MLST**     Multi-locus sequence typing

**mrca**     Most recent common ancestor

**Mya**     Million years ago

**Nal**     Nalidixic acid

**NICED**     National Institute for Cholera and Enteric Diseases, Kolkata, India

**NTS**     Non-typhoidal salmonellosis

**OUCRU**     Oxford University Clinical Research Unit, Hospital for Tropical Diseases, Ho Chi Minh City, Vietnam

**PFGE**     Pulsed-field gel electrophoresis

**PSU**     Pathogen Sequencing Unit at the Wellcome Trust Sanger Institute, Cambridge, UK

**SNP**     Single nucleotide polymorphism

**SPI**     *Salmonella* Pathogenicity Island

**Tn**     Transposon

**TTSS**     Type III secretion system

**VNTR**     Variable number tandem repeat