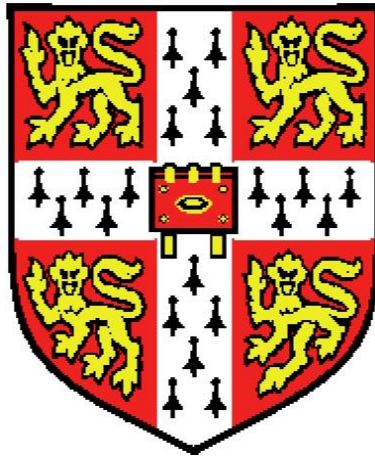# *De novo* assembly of the *var* multi-gene family in clinical samples of *Plasmodium falciparum*

Samuel Ayalew Assefa

Wolfson College

University of Cambridge

This dissertation is submitted for the degree of

*Doctor of Philosophy*

February 2013

# 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 stated here and in the text.

Sequence data used in this thesis was generated at the Sanger Institute by Research and Development and Sequencing production teams.

None of the work presented 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.

Samuel A. Assefa

February 2013

*To my family*

# Acknowledgements

# Abstract

In the malaria parasite *Plasmodium falciparum*, *Pf*EMP1 (*Plasmodium falciparum* erythrocyte membrane protein 1) is a protein that is exported to the surface of infected human red blood cells and encoded by ~60 *var* genes. *Pf*EMP1 plays a crucial role in parasite virulence and pathogenesis. It is also a target of host protective antibody responses that are avoided by the parasite by transcriptional switches between members of the *var* gene family resulting in antigenic variation of the surface expressed *Pf*EMP1.Thousands of malaria patient samples are being sequenced at the Sanger Institute's Malaria Programme to identify common polymorphisms but paradoxically some of the most variable sequences, such as *var* genes, are intractable due to high levels of polymorphism. Our understanding of *var* diversity in natural populations is thus limited to the DBL$\alpha$ domain, a conserved 300-400 bp region found in the majority of *var* genes studied so far. This thesis describes novel approaches developed to assemble full-length *var* genes from short reads of the Illumina sequencing platform.

The first part details an evaluation of existing assembly approaches through a comparative assessment of representative assembly tools. The results suggest that assembly of *var* genes in clinical samples using current methods is not practicaldue to a combination of factors including inherent sequence features (eg. high A+T content, low complexity, repeats and duplicates) and technical issues that affect quality of the raw sequence (eg. sequencing errors and uneven coverage). An alternative assembly strategy based on conserved sequence motifs was developed to address limitations of existing methods.

The second part investigates applications of short read sequencing to understand mechanisms of *var* gene diversity. Analysis of sequences from five progeny of the first genetic cross in *Plasmodium falciparum* between clones 3D7 and HB3 revealed evidence of ectopic-recombination as a mechanism for *var* gene diversity.

In the third and final part of the thesis, the iterative assembly approach developed in the first part is applied to a global collection of ~800 clinical isolates resulting in the first and largest collection of full-length sequences for ~50,000 *var*-contigs. Assembly results of *var* genes from these clinical samples were shown to have a higher repertoire-completeness (i.e. the number of contigs identified as *var* genes was close to the expected number of *var* genes), and contiguity (i.e. contig N50 size, largest contig size and open reading frame sizes were comparable with the expected values from previously completed genomes such as 3D7). Such availability of full-length *var* genes is a major progress towards understanding the population structure and diversity of *var* genes in natural populations with

Preliminarily analysis of *var*-contigs based on nucleotide and amino acid similarities (Chapter 5) revealed distinct clusters of highly conserved *var*-contigs within and between populations with percent-identities of up to 100% over their full length (i.e a match length of ~5-10 kb). The validity of these continent-transcending *var*-contigs was confirmed by looking at the sizes of open reading frames and aligning short reads back to *var*-contigs. Potential reasons for such continent-transcending *var*-contigs are explored in Chapter 6. These observations were surprising and potentially interesting as the majority of continent-transcending *var*-contigs were members of a group of *var* genes that are known to be associated with severe malaria.

# Contents

# List of Figures

# List of Tables

# Abbreviations

aa           amino acid

ACT          Artemis Comparision Tool

AT (A+T)     Adenine and Thiamine

BAM          Binary Alignment and Map

bp           base pair

CTV          Continent-transcending *var*-contig

DBL          Duffy binding-like

DNA          Deoxyribonucleic acid

DSB          double strand break

DSBR         double-strand break repair

G+C          Guanine and Cytocine

GA           Genome analyser

kb           Killo base

Mb           Mega base

nt           nucleotide

ORF          Open reading frame

PCR             Polymerase chain reaction

PfEMP1          *Plasmodium falciparum* erythrocyte membrane protein 1

RAM             Random Access Memory

SAM             Sequence alignment and Map

SDSA            Synthesis Dependent Strand Annealing

SNPs            Single nucleotide polymorphisms

TAR             transformation-associated recombination