

Causes and consequences of adult sepsis in Blantyre, Malawi

Thesis submitted in accordance with the requirements of the Liverpool School of Tropical Medicine for the degree of Doctor in Philosophy by Joseph Michael Lewis

August 2019

Contents

Preface	11
1 Introduction	13
1.1 Chapter Overview	15
1.2 Sepsis in sub-Saharan Africa	15
1.3 ESBL-E in sub-Saharan Africa	15
1.4 Conclusions	15
1.5 Thesis overview	15
1.6 Appendix	15
1.7 References	15
2 Methods	17
2.1 Chapter Overview	19
2.2 Study site	19
2.3 Clinical Study	19
2.4 Diagnostic Laboratory Procedures	19
2.5 Molecular methods	19
2.6 Bioinformatics	19
2.7 Statistical Analysis	19
2.8 Study Team	19
2.9 Data Collection and Storage	19
2.10 Ethical Approval, Consent and Participant Remuneration	19
3 A clinical and microbiological description of sepsis in Blantyre, Malawi	21
3.1 Chapter overview	22
3.2 Introduction and chapter aims	22
3.3 Methods	22
3.4 Results	22
3.5 Discussion	22

3.6 Conclusions and further work	22
4 Modelling to identify determinants of sepsis mortality	23
4.1 Chapter overview	24
4.2 Introduction and chapter aims	24
4.3 Methods	24
4.4 Results	24
4.5 Discussion	24
4.6 Conclusions and further work	24
4.7 Appendix	24
5 ESBL-E carriage in Malawian adults in health and disease	25
5.1 Chapter Overview	26
5.2 Introduction and chapter aims	26
5.3 Methods	26
5.4 Results	26
5.5 Discussion	26
5.6 Conclusions and further work	26
6 Whole genome sequencing of ESBL <i>E. coli</i> carriage isolates	27
6.1 Chapter overview	29
6.2 Introduction and chapter aims	29
6.3 Methods	29
6.4 Results	29
6.5 Discussion	29
6.6 Appendix	29
7 Using whole genome sequencing as a high-resolution typing tool to track ESBL-E carriage	31
7.1 Chapter overview	31
7.2 Introduction and chapter aims	31
7.3 Methods	31
7.4 Results	33
7.5 Discussion	41
7.6 Conclusions and further work	42
8 Longitudinal models of ESBL-E carriage	47
8.1 Chapter Overview	49
8.2 Introduction and chapter aims	49

CONTENTS	5
8.3 Methods	49
8.4 Results	49
8.5 Discussion	49
8.6 Conclusion and further work	49
8.7 Appendix	49
References	51

List of Tables

List of Figures

7.1	Core gene hierBAPS clusters	35
7.2	ESBL-clusters mapped back to phylogeny	37
7.3	Healthcare association of isolates mapped to phylogeny	38
7.4	Within- and between pairwise SNP distance of isolates	39
7.5	Within-patient ESBL-cluster and hierBAPS clustering as a function of time .	40
7.6	Summary statistics for ESBL-contig clusters	44
7.7	AMR genes, IS and plasmid replicon content of ESBL-clusters, part 1	45
7.8	AMR genes, IS and plasmid replicon content of ESBL-clusters, part 2	46

Preface

Placeholder

Chapter 1

Introduction

Placeholder

1.1 Chapter Overview

1.2 Sepsis in sub-Saharan Africa

1.2.1 Search strategy

1.2.2 Defining sepsis

1.2.3 Applicability of sepsis-3 definitions in sub-Saharan Africa

1.2.4 Sepsis epidemiology in sub-Saharan Africa

1.2.4.1 Incidence

1.2.4.2 Risk factors: the sepsis population in sub-Saharan Africa

1.2.4.3 Outcomes

1.2.5 Sepsis aetiology in sub-Saharan Africa

1.2.5.1 Bacterial zoonoses, Rickettsioses and arboviruses

1.2.5.2 HIV opportunistic infections: PCP, histoplasmosis and cryptococcal disease

1.2.6 Sepsis management

1.2.6.1 Early goal directed therapy

1.2.6.2 Evidence to guide antimicrobial therapy in sSA

1.2.6.3 Evidence to guide intravenous fluid therapy in sub-Saharan Africa

1.3 ESBL-E in sub-Saharan Africa

1.3.1 Search strategy

1.3.2 Introduction: definition and classification of ESBL-E

1.3.3 Global molecular epidemiology of ESBL-E: an overview

1.3.3.1 1980s-1990s: First identification of ESBL in nosocomial pathogens

1.3.3.2 1990s-2010s: Emergence and globalisation of CTX-M

Chapter 2

Methods

Placeholder

2.1 Chapter Overview

2.2 Study site

2.2.1 Malawi

2.2.2 Queen Elizabeth Central Hospital

2.2.3 Participating Laboratories

2.2.3.1 Malawi-Liverpool-Wellcome Clinical Research Programme

2.2.3.2 Malawi College of Medicine Tuberculosis Laboratory

2.2.3.3 Wellcome Trust Sanger Institute

2.3 Clinical Study

2.3.1 Entry Criteria

2.3.2 Study Visits and Patient Sampling

2.3.2.1 Enrollment assessment and first six hours

2.3.2.2 Subsequent visits

2.3.2.3 Blood, urine, and stool, sputum and CSF collection

2.3.2.4 Imaging: chest x-ray and ultrasound scanning

2.3.3 Outcomes and sample size calculations

2.4 Diagnostic Laboratory Procedures

2.4.1 Point of care diagnostics

2.4.2 Laboratory diagnostics

2.4.2.1 Haematology and biochemistry

2.4.2.2 Aerobic blood and CSF culture

2.4.2.3 Mycobacterial blood culture

Chapter 3

A clinical and microbiological description of sepsis in Blantyre, Malawi

Placeholder

3.1 Chapter overview

3.2 Introduction and chapter aims

3.3 Methods

3.4 Results

3.4.1 Study population

3.4.2 Baseline characteristics

3.4.3 Admission physiology and laboratory investigations

3.4.4 Aetiology

3.4.5 Treatment

3.4.6 Outcome

3.4.7 Determinants of mortality

3.5 Discussion

3.5.1 Demographics and outcome: significant longer-term mortality

3.5.2 Aetiology: TB dominates as a cause of sepsis

3.5.3 Determinants of mortality

3.5.4 Limitations

3.6 Conclusions and further work

Chapter 4

Modelling to identify determinants of sepsis mortality

Placeholder

4.1 Chapter overview

4.2 Introduction and chapter aims

4.3 Methods

4.4 Results

4.4.1 Exploring time-to antibiotics and IV fluid as determinants of mortality

4.4.2 Propensity score matching and subgroup analysis

4.5 Discussion

4.5.1 Limitations

4.6 Conclusions and further work

4.7 Appendix

Chapter 5

ESBL-E carriage in Malawian adults in health and disease

Placeholder

5.1 Chapter Overview

5.2 Introduction and chapter aims

5.3 Methods

5.4 Results

5.4.1 Study population

5.4.2 Exposures during the study period

5.4.3 ESBL-E colonisation

5.4.4 Associations of ESBL colonisation

5.5 Discussion

5.5.1 Limitations

5.6 Conclusions and further work

Chapter 6

Whole genome sequencing of ESBL *E. coli* carriage isolates

Placeholder

6.1 Chapter overview

6.2 Introduction and chapter aims

6.3 Methods

6.3.1 Bioinformatic pipeline

6.3.2 Global *E. coli* collection

6.3.3 Statistical analysis

6.4 Results

6.4.1 Samples and quality control

6.4.2 Phylogroup, MLST and core genome phylogeny of study isolates

6.4.3 Study isolates in a global context

6.4.4 Antimicrobial resistance determinants

6.4.4.1 β -lactam resistance

6.4.4.2 Quinolone resistance

6.4.4.3 Aminoglycoside resistance

6.4.4.4 Chloramphenicol, co-trimoxazole, tetracycline and other resistance determinants

6.4.4.5 Clustering and lineage association of AMR determinants

6.4.5 Plasmid replicons

6.5 Discussion

6.5.1 Genomic landscape of ESBL *E. coli* in Malawi: global diversity and high-risk clones

6.5.2 Antimicrobial resistance determinants: domination of *bla_{CTXM-15}* and emergence of carbapenemases

6.5.3 Study limitations

Chapter 7

Using whole genome sequencing as a high-resolution typing tool to track ESBL-E carriage

7.1 Chapter overview

7.2 Introduction and chapter aims

7.3 Methods

The rhierbaps v1.1.0 package in R[1] was used to cluster the core genome pseudosequence into sequence clusters (SCs). Two levels were used and these level 2 clusters used to test associations (see statistical analysis, below). To track putative mobile genetic elements ESBL-gene containing contigs were identified using BLASTn v2.7.0[2] of all contigs against the SRST2 database and then contigs containing any given ESBL gene were grouped by the ESBL gene they contained (for example, all *bla_{CTXM-15}* gene-containing clusters were grouped together), and each group clustered using cd-hit v4.6[3] to produce mutually exclusive ESBL-gene-containing contig clusters for each identified ESBL gene. Henceforth, these clusters will be referred to as ESBL-clusters, for brevity. In order to attempt to determine the biological significance of the identified ESBL-clusters (i.e. what kind of MGE element they are likely to represent), basic statistics were plotted (number of samples contained withing each cluster, length of longest contig in cluster in kbases, length distribution of all contigs is cluster relative to longest contig and distribution of sequence identity compared to the longest contig in the

cluster). Presence of insertion sequences (i.e compound transposons), AMR determinants and plasmid replicons were identified by using BLAST with default settings of each ESBL-cluster representative sequence (as determined by cd-hit i.e one, the longest, for each ESBL-cluster) against the insertion sequence finder (ISfinder) database and the SRST2 database, filtering such that sequence identify was greater or equal to 95%, taking the top hit (as determined by bitscore) for any given location if there were two overlapping hits, and visualising the results in ggenes v0.3.2. To assess lineage association, the ESBL-clusters were mapped back to the core gene phylogeny.

To explore hospital or community associations of any given *E. coli* clade, the location of isolation was first plotted against the phylogenetic tree; location of isolation was classified as hospital, community, or recent hospital discharge (defined as a date of isolation within 2 weeks of hospital discharge). This latter category was used because it is possible that a patient could acquire an ESBL-E clone in hospital but only be sampled once leaving hospital; using only hospital isolated and community isolated categories could therefore introduce bias. Hospital or community association of each sequence cluster was assessed using a Fisher's exact test of proportion of hospital associated samples (defined as sum of hospital isolated and recent hospital discharge) for the given sequence cluster as compared to proportion of hospital associated samples in the remainder of the samples, with a Bonferroni correction for multiple comparisons. $p < 0.05$ was again considered statistically significant.

To compare within-patient to between-patient conservation of bacteria (as represented by core genome alignment and sequence cluster) and ESBL-containing MGE (as represented by the ESBL-clusters) several approaches were taken. Firstly, I assessed whether either sequence cluster or ESBL-cluster were conserved within an individual at all. I hypothesised that any within-patient correlation is likely to be a function of time: samples closer together in time may be more likely to be similar. To assess if this was the case for bacteria, pairwise core genome pseudosequence SNP distance was calculated using snp-dists v0.4 (<https://github.com/tseemann/snp-dists>) for all samples and plotted against the time difference (in days) between samples, within and between patients, and with a smoothed curve fitted using a general additive model (GAM) with cubic splines. Because of significant overplotting, this was also plotted as a 2D density plot. Based on these plots, the within and between patient SNP distances were compared in two post-hoc defined groups binned by time distance between the samples (50 days or less vs. more than 50 days, these cutoffs determined from inspection of the pairwsie SNP distance vs time plots), and distributions compared with Kruskal-Wallace tests.

I then compared the within patient temporal clustering of ESBL-clusters and sequence clusters, by estimating the proportion of within-patient samples that contain the same ESBL-cluster

or sequence cluster, as a function of time; essentially a temporal auto-correlation function. To estimate this, I considered pairwise comparison of all within-patient samples. For any given time between samples, t I defined a window of $+/-5$ days and estimated the probabilities as the number of all within-patient sample pairs in the window $[t - 5, t + 5]$ that contained the same sequence cluster or ESBL-cluster divided by the total number of all within-patient sample pairs within that time window. Exact binomial confidence intervals for these proportions were generated and probabilities plotted as a function of time. In order to estimate the probability of two samples containing the same sequence cluster or contig-cluster purely by chance, 1000 sample pairs were randomly drawn from all samples with replacement and the proportion of these samples that contained the same sequence cluster or ESBL-cluster calculated.

Finally, to inform the question as to what the likely unit of transmission in this system is, I assessed what was most conserved within patients, in pairwise sample comparison: bacteria (as represented by core gene sequence cluster), ESBL-containing MGE (as represented by ESBL-cluster), or both. Simple proportions in all-against-all pairwise comparison - stratified by whether between-patient or within-patient - were calculated: the proportion of samples that contain the same core gene sequence cluster only, the proportion of samples contain the same ESBL-cluster only, and the proportion that contain both sequence cluster and ESBL-cluster. Proportions were compared between within and between-patient strata in these three groups using Fisher's exact test, with $p < 0.05$ considered statistically significant.

7.4 Results

7.4.1 Testing metadata associations: SNP distance, hierBAPS sequence clusters and ESBL-clusters

Finally, in order to test metadata associations of bacterial lineages or MGE, I used several techniques: considering core gene SNP distance between isolates to infer continuous carriage and/or transmission events, and clustering core gene pseudosequences and ESBL-containing contigs into mutually exclusive groups which can then be used to test associations. Below, I first describe the outcomes of the clustering algorithms used, before describing tests of association with metadata.

7.4.1.1 Hierarchical BAPS clustering of core gene pseudosequences

The hierarchical BAPS algorithm clustered the core gene alignments into 15 level one (top level) clusters, denoted sequence clusters A-O, and a total of 48 level two (lower level)

clusters, denoted sequence clusters 1-48 that were almost exclusively monophyletic and often corresponded closely to the multilocus sequence types (STs, Figure 7.1A). Intracluster pairwise SNP distance varied (Figure 7.1B) but the clusters were often reasonably clonal: SC6, SC8 and SC23, for example (the three largest clusters) had median (IQR) intragroup pairwise SNP distance of 62 (34-97), 326 (18-378) and 18 (11-24) respectively.

7.4.1.2 ESBL-clusters

The 473 samples contained 486 ESBL genes (Figure 7.2A); 5 genes only occurred once in the collection and so no attempt was made to cluster them. Of the remaining 481 genes pairs, BLAST failed to identify the ESBL-gene containing contig in 2 samples (one in which *ARIBA* had identified *bla-CTXM-15* one *bla_{CTXM-27}*), but identified the remaining 479 ESBL genes on 478 contigs, with perfect agreement with *ARIBA* as to which AMR gene was present in which sample. Only one contig carried two ESBL genes: *bla_{CTXM-3}* and *bla_{CTXM-15}*; the remaining 477 contigs contained one. The *cd-hit* algorithm grouped the 477 unique contigs into 99 clusters (Figure 7.2B). In total, over 90% of the ESBL-genes (432/479 [90%]) were contained in the 52 largest contig clusters.

The *cd-hit* algorithm selects one member of a cluster (the longest) as the representative. The structure of these representative contigs was explored in an attempt to understand type of MGE they were likely to represent. The length of the representative clusters was very variable, ranging from 1.8kbp to 905.8kbp, with median (IQR) 46.1kbp (11.1-215.5kbp). The other cluster members were usually fragments of these representative contigs with varying sizes - a median (IQR) 60% (36-100%) of the representative contig length - but had high sequence identity, median (IQR) 100.0% (99.7-100.0%) (Figure 7.6 in the appendix to this chapter).

I then explored the insertion sequence (IS), AMR gene and plasmid replicon content of the representative contig for each cluster using BLAST against the SRST2, ISfinder and Plasmidfinder databases (Figures 7.7 and 7.8 in the appendix to this chapter). Every ESBL gene was closely associated with at least one IS, commonly ISEcp1, IS26 and IS903B. IS26 was frequently associated with an apparent 108bp fragment of a *catB4* chloramphenicol resistance determinant. Some ESBL-genes were associated with particular IS; *bla_{CTXM-15}*, *bla_{CTXM-9}* and *bla_{CTXM-1}*, for example were very commonly associated with ISEcp1, whereas *bla_{SHV-12}* was associated with IS26. ESBL genes were not infrequently associated with other resistance determinants, including commonly *bla_{CTXM-15}* with *bla_{TEM-95}*. Plasmid replicons were occasionally identified, including an IncFIB plasmid carrying *bla_{CTXM-15}* and an IncQ1 plasmid carrying *bla_{CTXM-27}*. It is clear that the same configuration of AMR genes and IS are seen across different contigs, despite a varying backbone, implying historical transposition events. Finally, to assess lineage associations of the identified ESBL-clusters, I mapped the

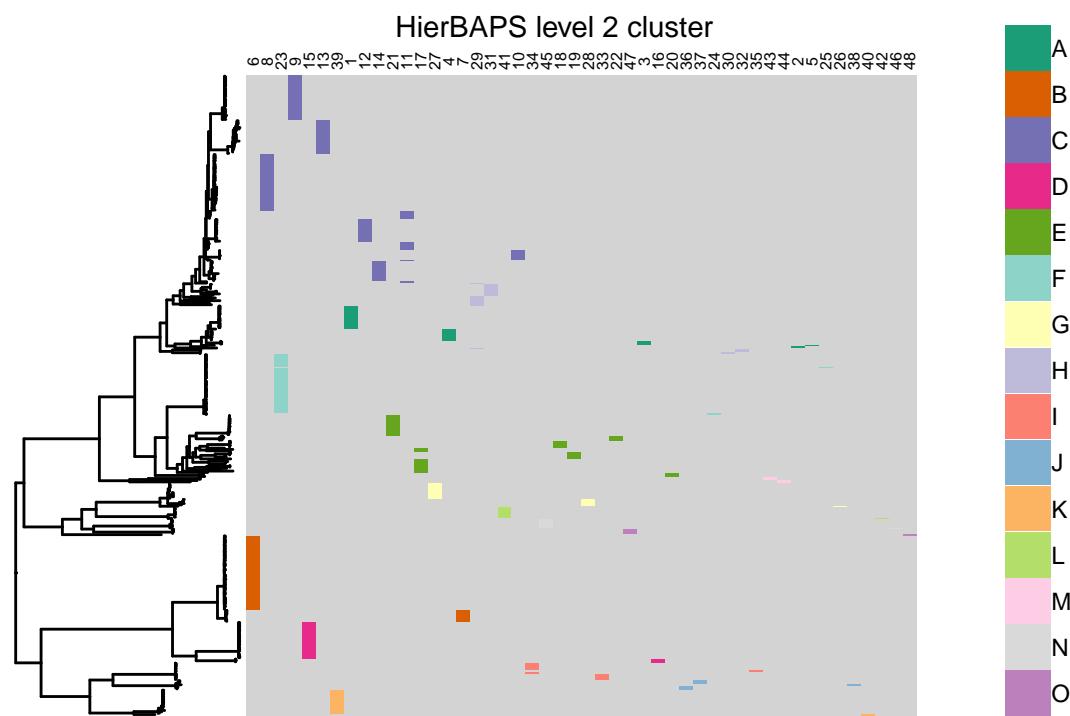
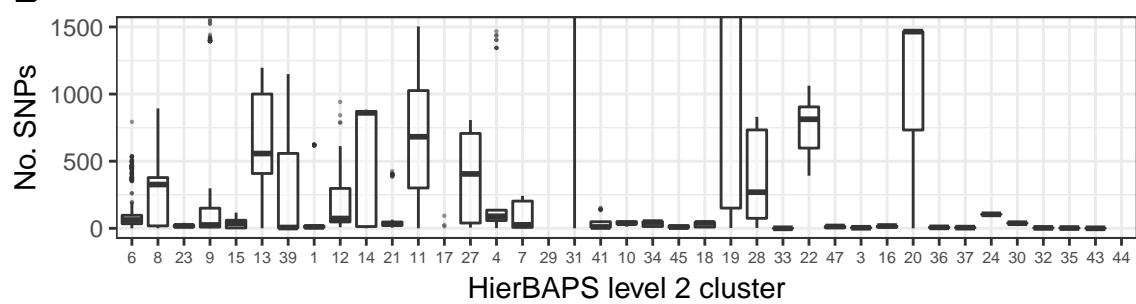
A**B**

Figure 7.1: A: Core gene hierarchical BAPS clusters mapped back to phylogeny. Heatmap shows level 2 (lower level) with colour denoting level 1 (top level) cluster membership. B: Intracluster pairwise SNP distance for level 2 sequence clusters. Axis restricted to 0-1500 SNPs and as result SC17 (median 6881 SNPs), SC29 (median 2970 SNPs), SC31 (median 2970 SNPs) and SC44 (median 12322 SNPs) boxes are not shown. Boxplots show median and IQR, whiskers show 1.5 times IQR, and outliers are points falling beyond whiskers.

clusters back to the tree, and found that there was a strong lineage association (Figure 7.2C).

7.4.1.3 Assessing for healthcare-associated lineages

Having clustered bacteria and MGE using *hierBAPS* and *cd-hit* respectively, I then mapped the location of sample collection back to the phylogeny and used the *hierBAPS* SCs to assess for healthcare associated lineages (Figure 7.3). In general, healthcare-associated isolates were distributed throughout the tree and across all SCs, rather than there being a clear hospital-associated lineage. The exception to this was SC23, corresponding to MLST 410, which was slightly more likely to be healthcare associated. When comparing the proportion of healthcare associated samples within each SC to the remained of samples, SC23 had a statistically significantly increased proportion of healthcare associated samples on Fisher's exact test ($p = 6.3 \times 10^{-4}$, threshold of significance following Bonferroni correction 1.0×10^{-3}), though it was by no means health-facility restricted: 50% (21/42) of SC23 samples were isolated in the community.

7.4.1.4 Assessing for within-patient conservation of lineage or MGE

To answer the question as to what elements (bacteria or MGE) are conserved within individuals across time I first compared all-against-all pairwise SNP distance between and within patients; first as a scatter plot, and then, because of significant overplotting, as a density plot (Figure 7.4). This suggested that there are a cluster of points close to the origin in the within-patient plot that are not seen in the between-patient plot: before approximately 50 days, there are more similar within-patient isolates than seen in the between-patient isolates. Dichotomising time at 50 days (based on inspection of the density plots) and performing a Kruskal-Wallace test found a statistically significant difference between the before 50 day and after 50 day pairwise SNP distance distribution in the within patient stratum ($p = 0.008$) but not in the between-patient stratum ($p = 0.07$). After 50 days, the distribution of between- and within-patient SNP distances are similar ($p = 0.45$). However it is clear from the plots that even at small t and within-participant, there is significant diversity in the SNP distances, and that some isolates close together in time, within the same participant, are only distantly related.

Having confirmed that there is a signal for within-participant temporal conservation of ESBL-E, I then sought to determine if the sequence clusters and ESBL-clusters were similarly conserved over time, and if so, which was the more conserved. The proportion of pairwise within-patient samples that contained the same ESBL-cluster and sequence cluster were significantly greater than would be expected by chance when the time between the samples is less than 35 days for sequence cluster and 32 days for ESBL-cluster (Figure 7.5A). After this time, the lower

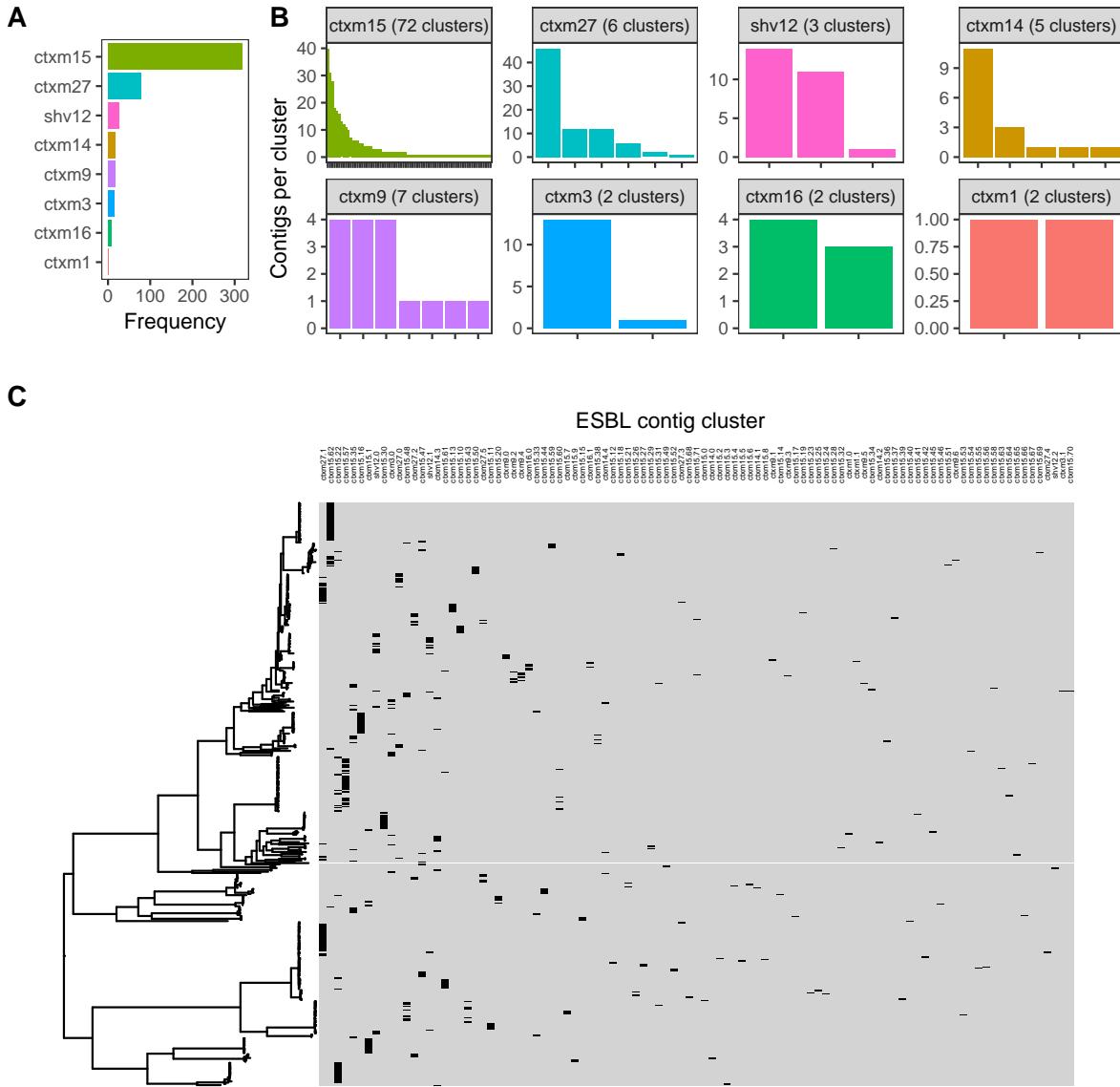


Figure 7.2: A: Frequency distribution of ESBL genes in included samples. B: Frequency distribution of samples per ESBL-cluster, stratified by gene. C: ESBL-cluster membership mapped back to phylogeny.

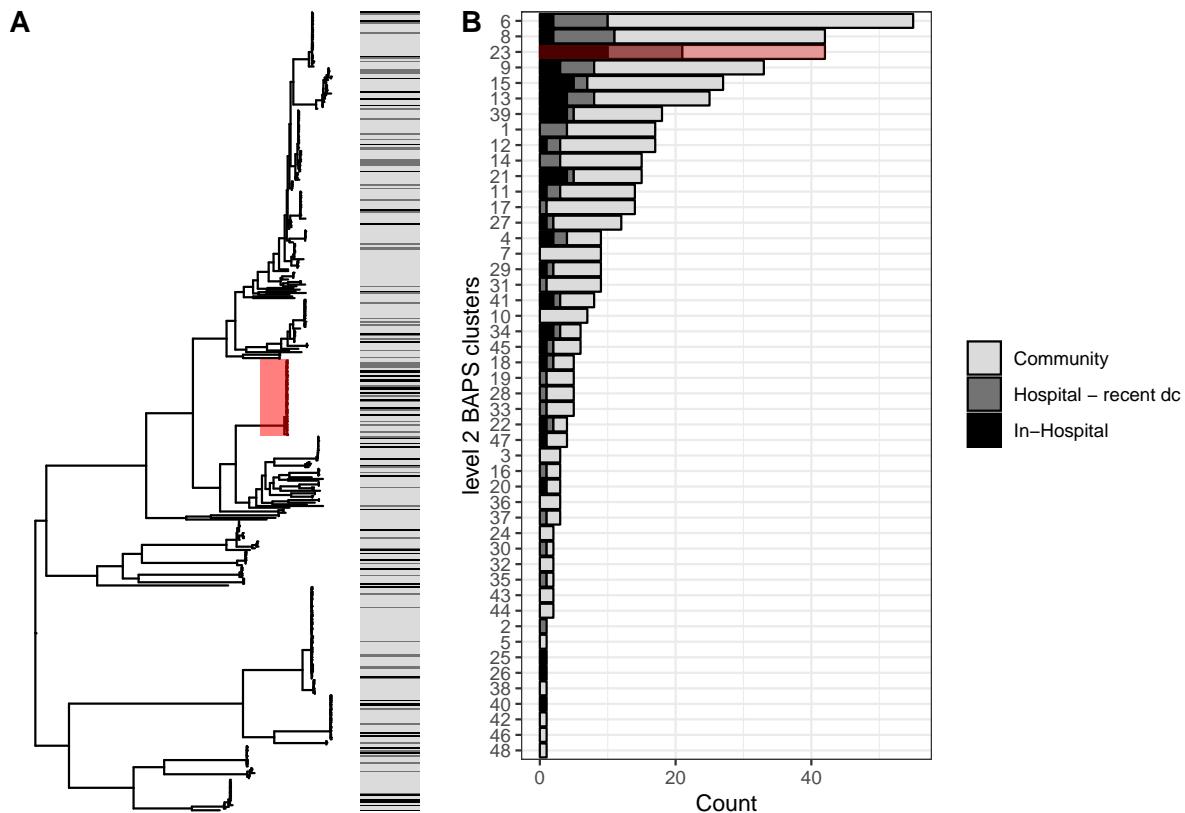


Figure 7.3: A: Location of sample isolation mapped back to phylogeny B: Distribution of location of sample isolation stratified by hierBAPS cluster. In each case, community isolates include those cultured from samples collected on the first day of hospital admission, in-hospital isolates are from patients who have been hospitalised > 24 hrs and recent discharge isolates are from patients who have been discharged from hospital within the last 2 weeks. Sequence cluster 23, highlighted in red, showed a statistically significant association with hospitalisation (see text).

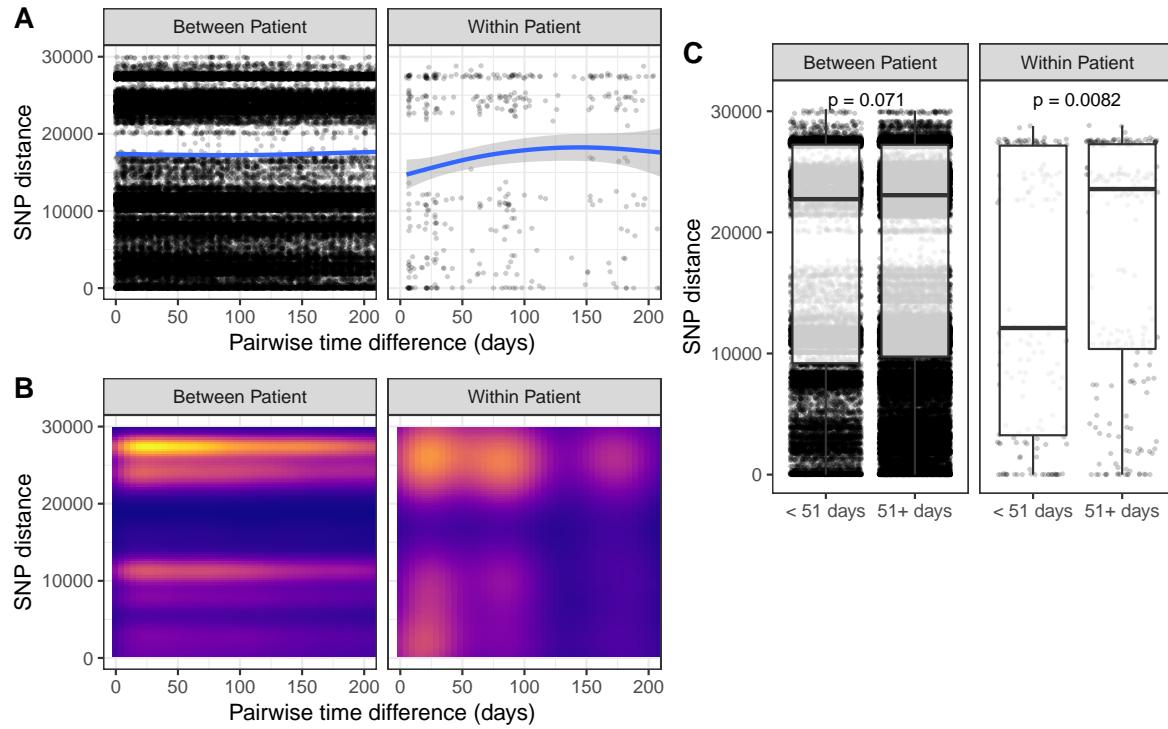


Figure 7.4: Within and between participant pairwise SNP distances. A: Scatterplot of pairwise SNP distances as a function of time with GAM model fitted curve. B: Pairwise SNP distance as function of time as a 2D density plot, showing cluster of isolates close to origin that are close together in time and SNP-distance. C: Pairwsise SNP distance distribution before and after 50 days, within and between patients, showing statistically significant descreas ein pairwise SNP distance within patients before 50 days. After 50 days, between- and within- patient distributuions are similar.

confidence interval of the sequence cluster and ESBL-cluster curve crossed the proportion of samples that would be expected to be the same by chance, suggesting that, after 35 or 32 days, the chance of any two within-patient samples having the same sequence cluster or ESBL-cluster (respectively) is the same as if the two samples were randomly picked from the data set without regard to patient. The two curves have a very similar appearance; to address the question of which element is most conserved within an individual - sequence cluster, ESBL-cluster, or both - I performed an all-against-all pairwise comparison of which elements were conserved (Figure (Figure 7.5C), and found that only ESBL-cluster and sequence cluster together are conserved within patients at a significantly greater proportion than between patients ($p = 1.1 \times 10^{-12}$).

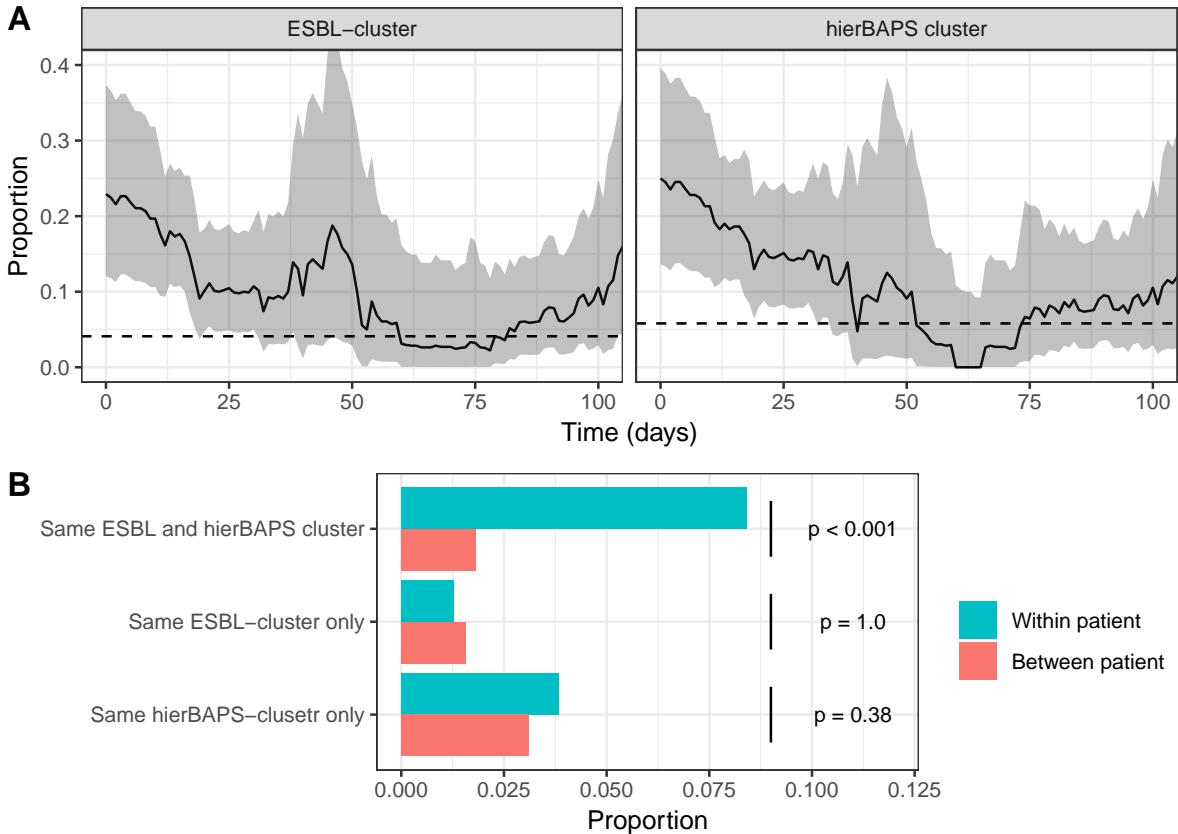


Figure 7.5: Probability of any two samples from within a given participant containing the same ESBL-cluster (A, left panel) or being a member of the same hierBAPS cluster (A, right panel). Time is windowed at ± 5 days around the time indicated on the x axis. Dotted line is the probability that two samples would belong to the same group by chance, constructed by randomly sampling 1000 sample pairs. B: proportion of samples that contain the same herBAPS cluster alone, or ESBL-cluster alone, or both, demonstrating that the ESBL cluster-hierBAPS cluster pairing is the most conserved of the three.

7.5 Discussion

7.5.1 Drivers of ESBL-E carriage: true acquisition versus enrichment

The diversity of healthcare associated isolates was largely contained within the diversity of community isolates, rather than apparent hospital acquisitions being restricted to a single clade or clone. The exception to this was SC23, which corresponded to ST410, and was more likely to be healthcare associated. This could be consistent with the hypothesis that it is a recently arrived high-risk clone, which may be, at least initially, hospital-associated. Even so, it is clearly not hospital restricted, with half of the ST410 isolates being isolated from the community. This result could be explained by two hypotheses: first, that these are true transmission events that are occurring within the hospital, and that the diversity of ESBL *E. coli* within the hospital is the same as the community; or, second, that these “hospital acquisitions” are actually minority variant *E. coli* present in the microbiota (and therefore acquired in the community) at a low abundance and hence not detected by culture, and enriched for with antimicrobial exposure in hospital. This latter hypothesis is perhaps supported by the fact that the antimicrobial-exposed group of patients in this study have an early and rapid increase in ESBL-E carriage prevalence that is not seen in the antimicrobial-unexposed group. Distinguishing between these two hypotheses is important as they would each require a different intervention: hospital infection control in the former case, or strategies to protect the microbiota from the deleterious effect of broad spectrum antimicrobials (such as pre- or probiotics, or oral β -lactamases) in the latter.

By forming sequence clusters and ESBL-clusters, I was able to demonstrate that both bacteria and MGE are conserved together, within-patient, over time. Some previous longitudinal studies of ESBL-E found that *E. coli* STs tended to vary over time but that in many cases ESBL gene and plasmid replicon remained the same, which could be due to a conserved MGE transferring between bacteria[4]. Given my findings, this is unlikely to be the case. Though not directly addressed in this study it is possible to speculate therefore that the unit of transmission of ESBL between patients is likely to be the bacterium rather than, for example, horizontal gene transfer of ESBL genes on plasmids or transposons. The within-patient correlation of SC and ESBL-cluster lasts only for 32-35 days; two samples from a single patient more than 35 days apart are as likely to contain the same SC/ESBL cluster as two samples from two different patients. This implies either an exogenous re-exposure or some other endogenous mechanism whereby the dominant ESBL strain is replaced by a minority variant from within the microbiota. Of note, the timescale of SC replacement - occurring after 35 days - is consistent with the mean time taken to revert to the ESBL-negative state from ESBL-positive from the longitudinal Markov models (26 [95% CI 12-58] days), perhaps

lending support to the re-exposure hypothesis. It is important to note also that even at a maximum (i.e. with samples that are days apart), only around 20% of within-patient samples contain the same SC/ESBL-cluster. This is many times more than would be expected by chance, but still implies that at any time point there is significant diversity of ESBL *E. coli* strains, that have been missed by only taking forward one colony pick for sequencing. This is consistent with studies that have performed multiple colony picks on stool samples enriched for ESBL-E, and have found widespread diversity of STs and ESBL genes[5].

7.5.2 Limitations

I selected a global collection of *E. coli* based on what was available but, as described above and in common with many analyses of this type, this is a biased collection. This must be borne in mind when interpreting the global phylogeny. There are inherent limitations in the short-read Illumina sequencing that was carried out: assembly of areas with multiple nucleotide repeats (as found in plasmids and transposable elements in particular) is difficult or impossible, making it impossible to fully characterise the MGE in this dataset upon which the AMR genes are carried. I have attempted to address this difficulty by defining ESBL-clusters as a proxy for MGE, but this is by its nature flawed. Some of the assembled contigs are short and likely represent transposons; the same transposons have likely inserted into multiple plasmids in the past and as such, these short contigs may cluster with other sequences that would be seen to be very different, were a full assembly available. In addition, the biological significance of these ESBL-clusters is not clear. It is not possible to say with certainty what they represent (e.g. plasmid) as they are only fragments. Nevertheless, the fact that I have seen within-patient associations of the ESBL-clusters lends some support to their use, as erroneous clustering would be expected to bias any associations towards to null.

7.6 Conclusions and further work

It is possible that apparent constant colonisation actually represents frequent re-exposure on the timescale of around 35 days, and that apparent hospital acquisitions, are, in fact, unmasking events due to enrichment of the microbiota for ESBL-E by antimicrobial exposure.

Many questions remain unanswered and further work is necessary. Shotgun metagenomic sequencing of stool would allow testing of the competing acquisition and unmasking hypotheses of rapid increase in ESBL-E prevalence by defining the microbiota and total AMR gene content pre-, during and post- antimicrobial exposure, as well as providing an opportunity to explore the role of the microbiota to colonisation resistance to ESBL-E. Long read sequencing would

allow a proper characterisation of the MGE that carry ESBL genes in the Malawian context, giving the resolution necessary to truly track MGE within and between patients and strains, as well as to address questions of expression of genes such as *catB4* that seem to correlate poorly with phenotype, by examining e.g. promoter regions. Short-read sequencing of the *Klebsiella pneumoniae* isolates from this study would allow a comparison between the mechanisms of AMR and MGE prevalent in this species as compared to *E. coli*, and assess the extent to which horizontal gene transfer between the two is driving ESBL spread in Blantyre. Collating publicly available ST167 and ST410 genomes would allow the construction of a phylogeny that would allow insight into the epidemiology of these putative recently arrived high risk clones. Finally, incorporating the resolution afforded by sequencing into the longitudinal Markov models may provide new insights in to the dynamics of ESBL-E carriage. This is the subject of the next chapter.

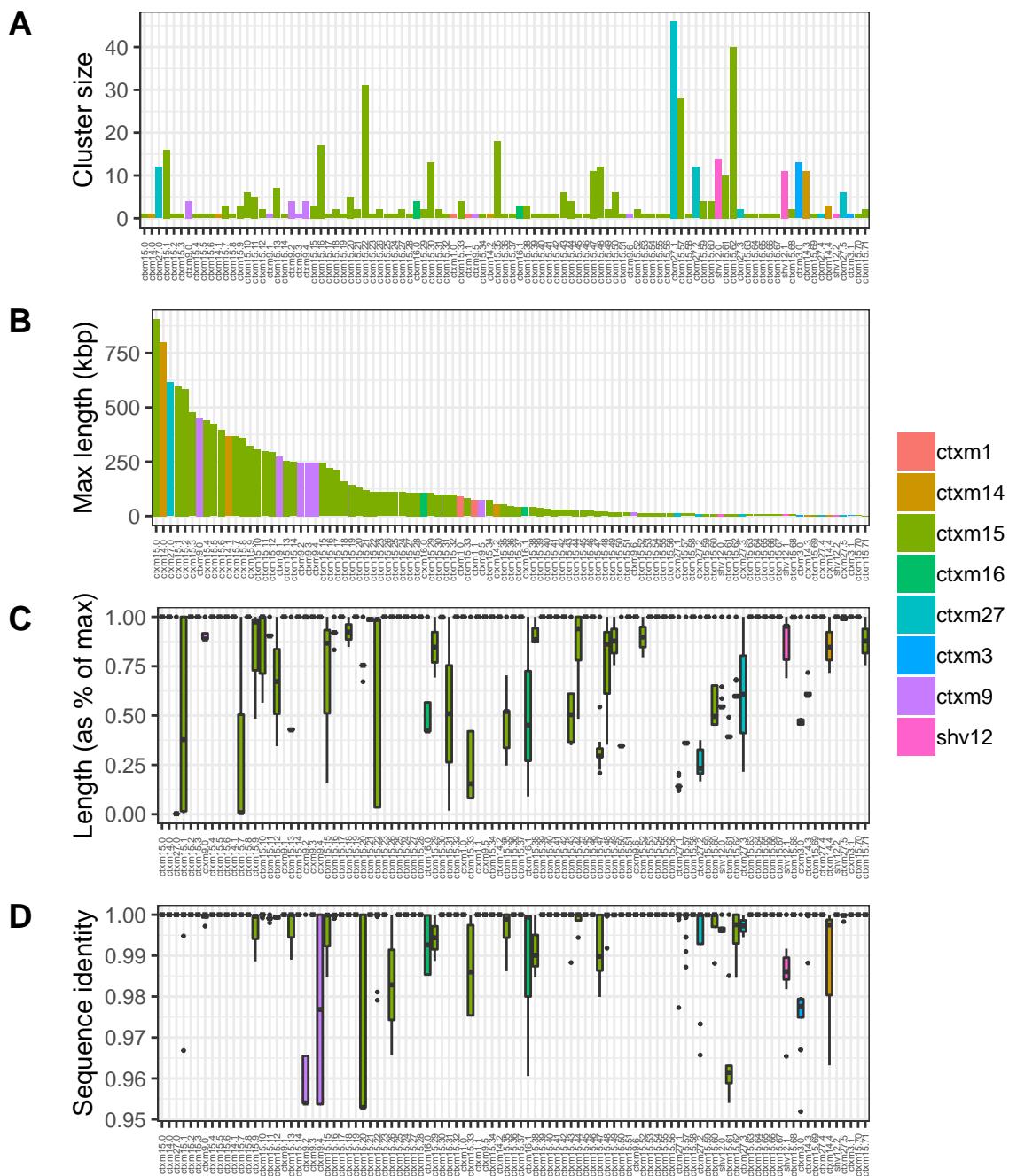


Figure 7.6: Summary statistics for 99 ESBL-containing contig clusters as determined by *cd-hit*. A: Number of contigs per cluster. B: Length (kbp) of longest sample in each cluster. This is defined as the cluster representative sample by *cd-hit* to which all other samples are compared for the purposes of length and sequence identity. C: Distribution of contig lengths by cluster expressed as a proportion of longest contig length. D: Distribution of sequence identity of cluster members compared to representative member, by cluster.



Figure 7.7: AMR genes, insertion sequences (IS) and plasmid replicons identified in the representative contig of each contig cluster, stratified by by ESBL gene and ordered by number of samples of cluster. IS26 is very frequently associated with a 108bp fragment of *catB4* chloramphenicol resistance gene, shown as a red fragnemt within the green IS26 element. A: *blaCTXM15*, B: *blaCTXM27* , C: *blaSHV12*. Plots show furthest IS/AMR gene or plasmid replicon up to +/- 10,000bp from the ESBL gene of interest.

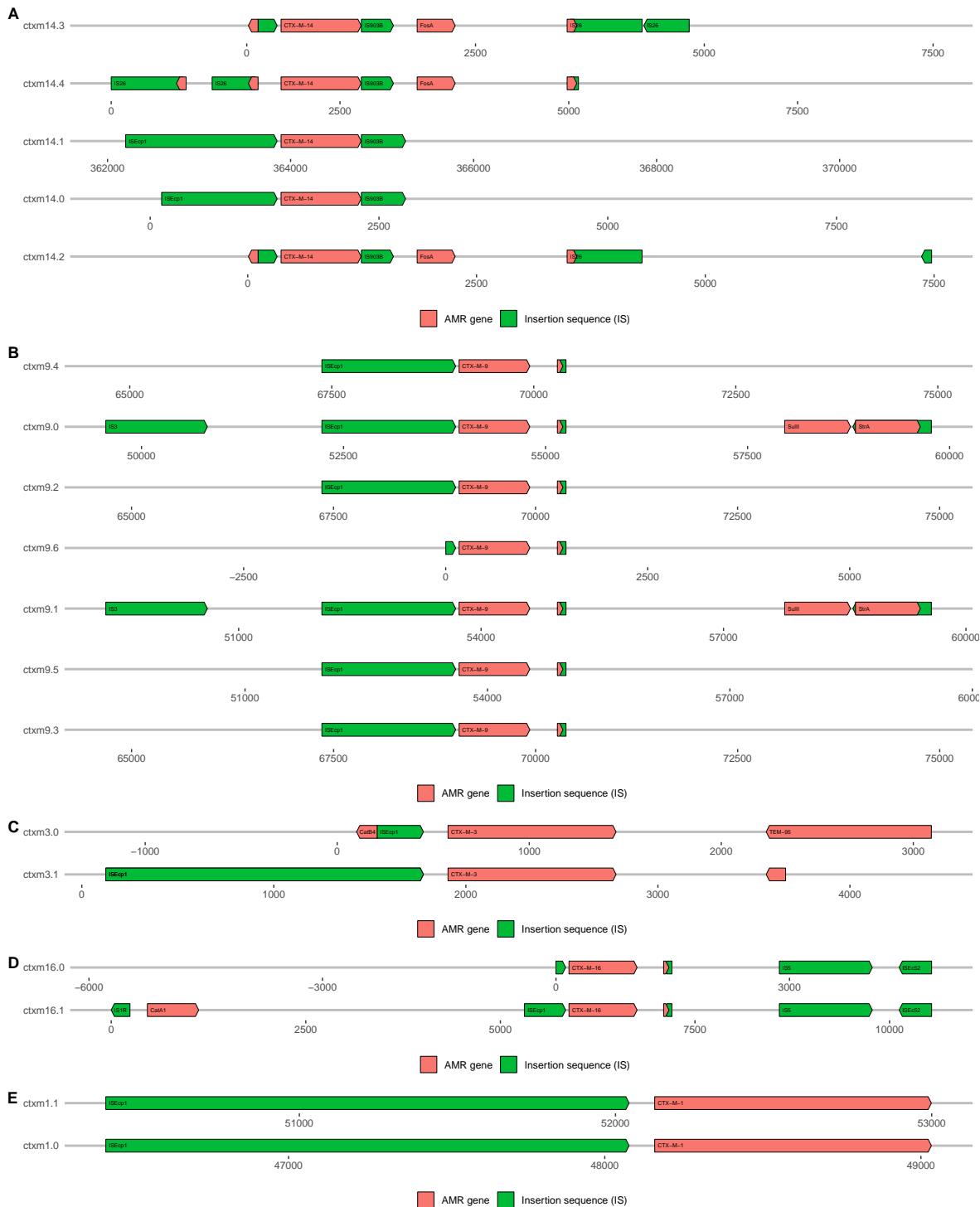


Figure 7.8: AMR genes, insertion sequences (IS) and plasmid replicons identified in the representative contig of each contig cluster, stratified by by ESBL gene and ordered by number of samples in cluster. IS26 is very frequently associated with a 108bp fragment of *catB4* chloramphenicol resistance gene, shown as a red fragnemt within the green IS26 element. A: *blaCTXM14*, B: *blaCTXM9* , C: *blaCTXM3*, D: *blaCTXM16*, E: *blaCTXM1*. Plots show furthest IS/AMR gene or plasmid replicon up to +/- 10,000bp from the ESBL gene of interest.

Chapter 8

Longitudinal models of ESBL-E carriage

Placeholder

8.1 Chapter Overview

8.2 Introduction and chapter aims

8.3 Methods

8.3.1 Developing the models used in this chapter

8.3.2 General form of likelihood

8.3.3 Markov model likelihood

8.3.4 Incorporating covariates: a proportional hazard model

8.3.5 Building and fitting models

8.3.6 Assessing goodness of fit

8.3.7 Exploring differences in carriage dynamics by bacterial species and *E. coli* genotype

8.3.8 Simulations from the posterior

8.4 Results

8.4.1 The effect of antibacterials and hospitalisation on ESBL-E carriage

8.4.2 Exploring bacterial species and genotype differences in carriage dynamics

8.4.3 Simulation of different antibacterial and hospitalisation scenarios

8.5 Discussion

8.5.1 Limitations

8.6 Conclusion and further work

8.7 Appendix

References

- 1 Cheng L, Connor TR, Siren J *et al.* Hierarchical and Spatially Explicit Clustering of DNA Sequences with BAPS Software. *Molecular Biology and Evolution* 2013;30:1224–8. doi:10.1093/molbev/mst028
- 2 Altschul SF, Gish W, Miller W *et al.* Basic local alignment search tool. *Journal of Molecular Biology* 1990;215:403–10. doi:10.1016/S0022-2836(05)80360-2
- 3 Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 2006;22:1658–9. doi:10.1093/bioinformatics/btl158
- 4 Duijkeren E van, Wielders CCH, Dierikx CM *et al.* Long-term Carriage of Extended-Spectrum β -Lactamase-Producing *Escherichia coli* and *Klebsiella pneumoniae* in the General Population in The Netherlands. *Clinical Infectious Diseases* 2018;66:1368–76. doi:10.1093/cid/cix1015
- 5 Stoesser N, Sheppard AE, Moore CE *et al.* Extensive Within-Host Diversity in Fe-cally Carried Extended-Spectrum-Beta-Lactamase-Producing *Escherichia coli* Isolates: Implications for Transmission Analyses. *Journal of clinical microbiology* 2015;53:2122–31. doi:10.1128/JCM.00378-15