



INTRODUCTION

Shiga toxin-producing *Escherichia coli* (STEC) are a group of zoonotic, foodborne pathogens defined by the presence of phage-encoded Shiga toxin genes (*stx*) [1]. STEC cause gastrointestinal disease in humans and symptoms include severe bloody diarrhoea, abdominal pain and nausea. In 5-15% of cases infection leads to Haemolytic Uremic Syndrome (HUS), characterised by kidney failure and/or cardiac and neurological complications [1].

STEC O157:H7 genomes range from 5.4Mbp to 5.6Mbp in size, and a high proportion (9-15%) is comprised of mobile genetic elements and prophages [2].

Due to the limitations of short-read sequencing technologies in handling the homologous regions of the STEC chromosome, information and context regarding inter- and intra-outbreak variation in prophages, structural variation and context surrounding plasmid content is lost. As a result, the prophage population of STEC O157 is largely unexplored.

In this project we used a combination of short-read Illumina and long-read Oxford Nanopore Technology (ONT) sequencing data to generate complete genome assemblies giving us access the accessory genome to characterise and compare the prophage content within domestic UK STEC O157:H7.

METHODS

DNA extraction was performed using a Qiagen Qiasymphony followed by library preparation using the Nextera XP kit followed by sequencing on the Illumina HiSeq 2500.

DNA extraction was also performed, using Revolugen's Fire Monkey kit followed by library preparation using SQK-RBK004 (Rapid) kit and sequencing on the Oxford Nanopore Technologies (ONT) MinION on a FLO-MIN106D flow cell.

Nanopore basecalling, read trimming and read filtering were performed using Guppy v3.2.10 FAST – Guppy v6.3.8 FAST, Porechop v0.2.4^[4] and Filtlong v2^[5] respectively.

Nanopore reads were assembled using Flye v2.9^[6] and the draft was corrected using four iterations of Racon v1.4.20^[7] (ONT reads), Medaka v0.1.0^[8] with an STEC-specific model (ONT reads), Pilon v1.23^[9] (Illumina reads) and finally Racon v1.4.20^[7] (Illumina reads).

Prophages sequences were collected manually from annotated finalised assemblies using Prokka v1.14.6^[10] and compared in a pairwise format using Mash v2.2.2^[11].



RESULTS

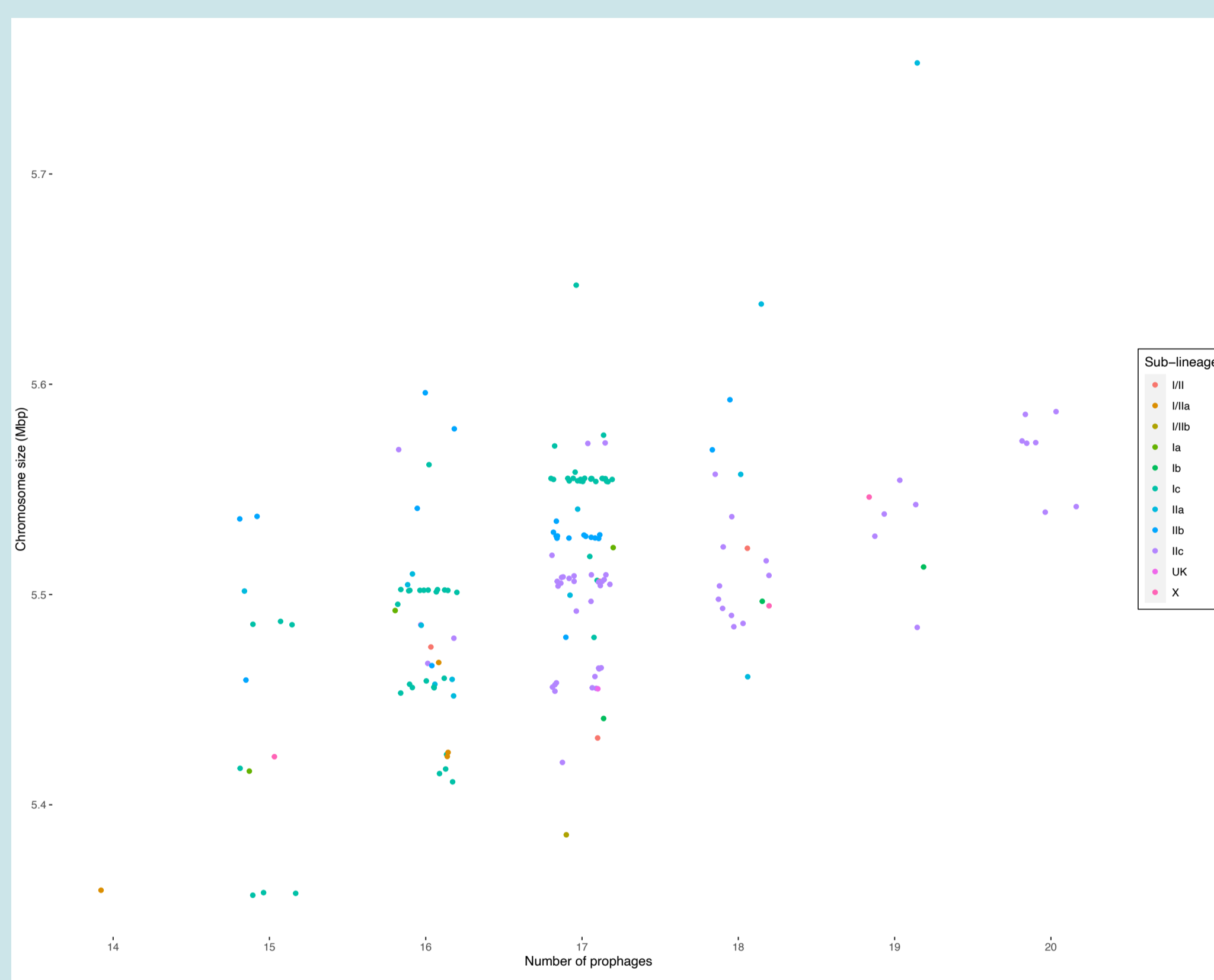
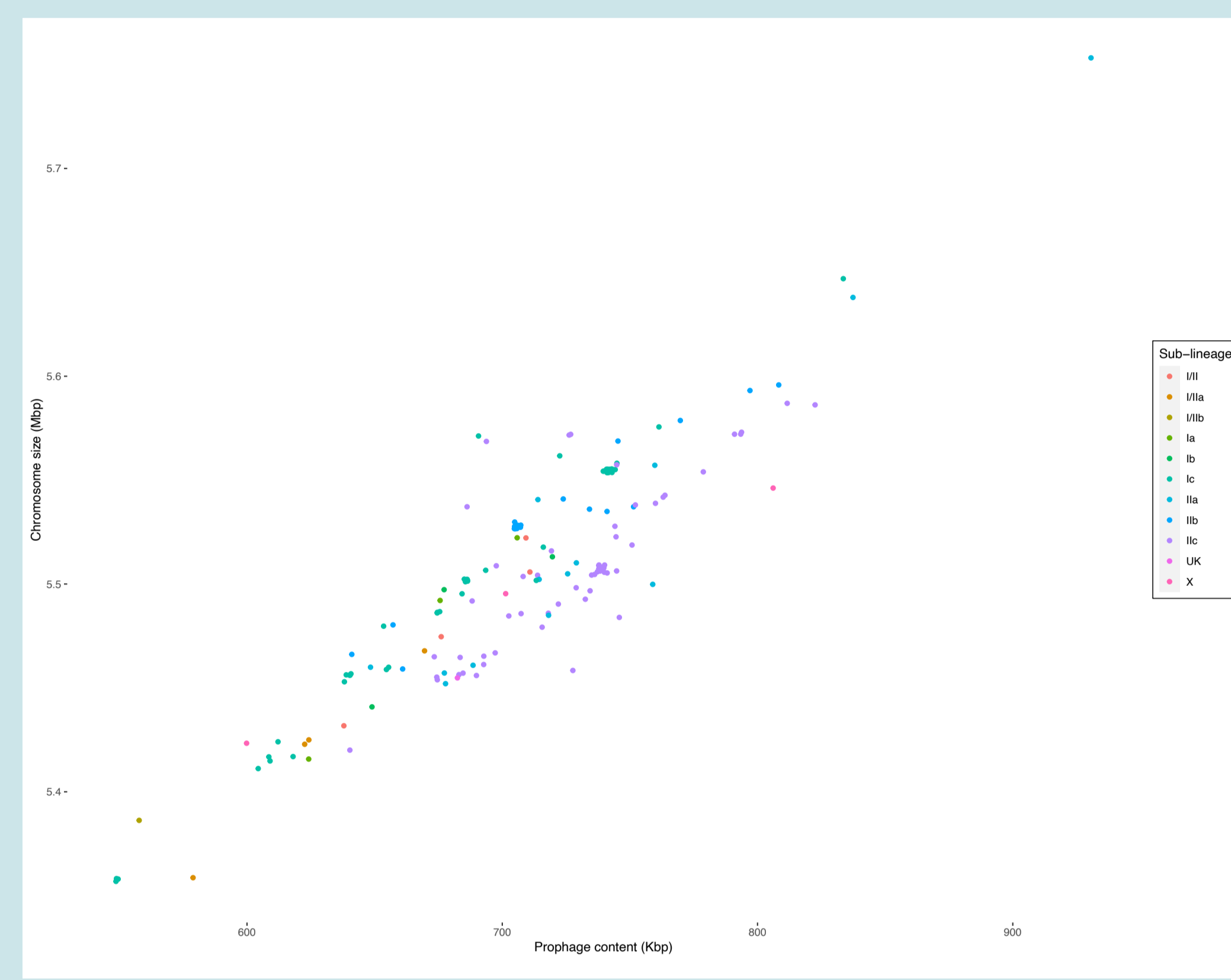


Figure 1. Left – Chromosome size vs number of prophages.



Right – Chromosome size vs total prophage content (kbp).

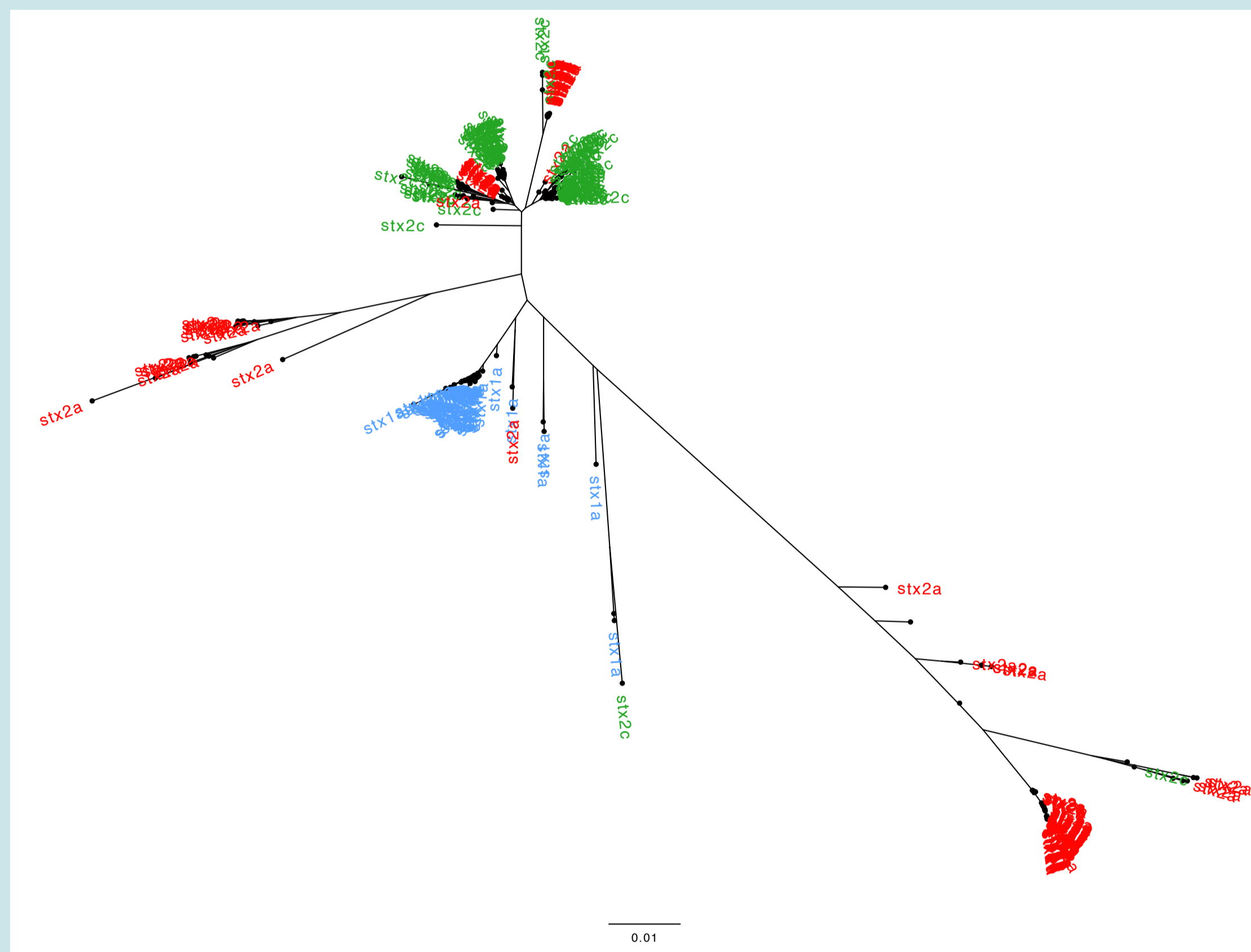


Figure 2. Neighbour joining tree based on Jaccard distances of all *stx*-encoding prophages ($n=308$) within the 176 STEC genomes sequenced. *stx2a*, *stx2c* and *stx1a* are coloured red, green and blue respectively.

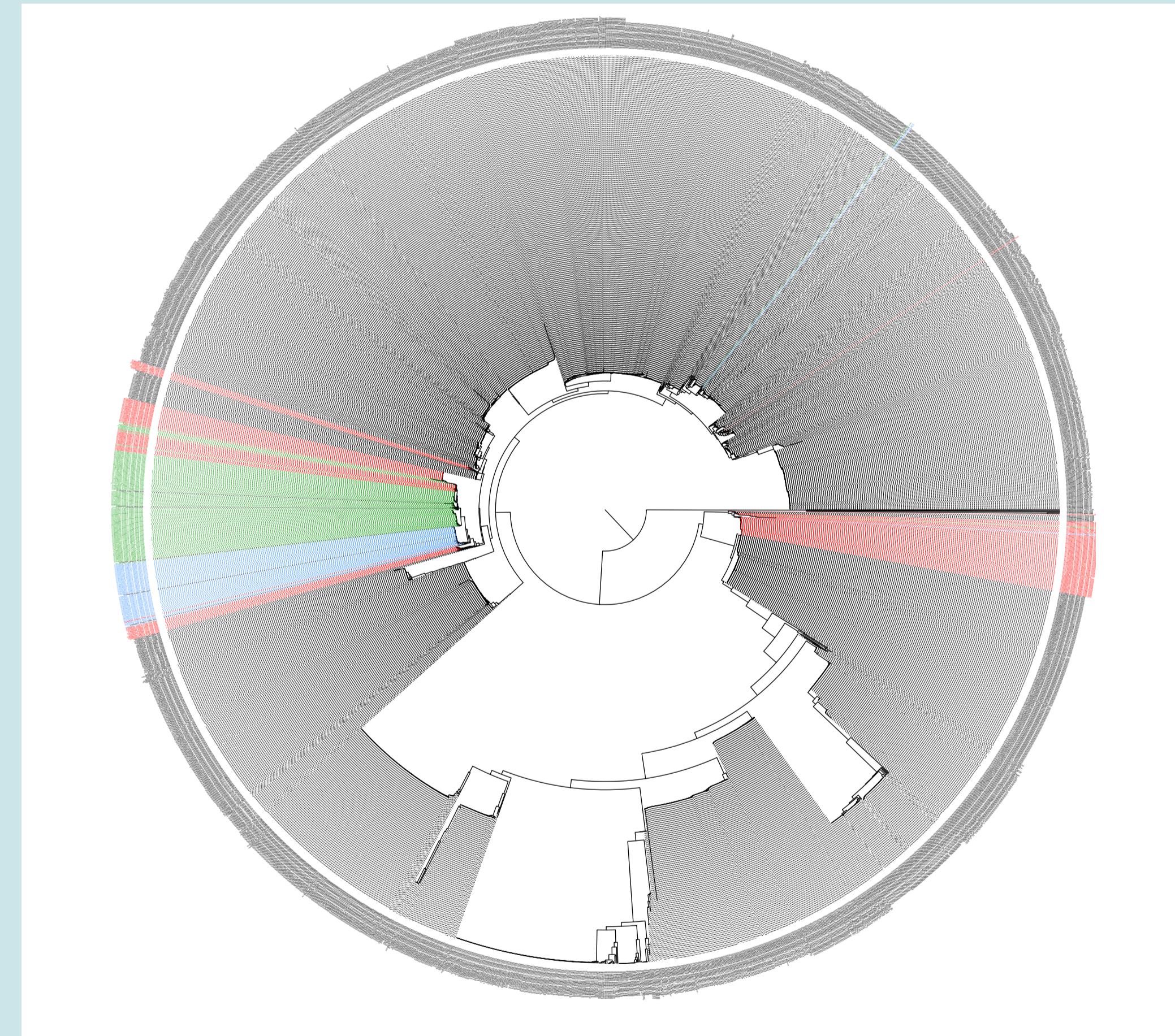


Figure 3. Neighbour joining tree based on Jaccard distances of all prophages ($n=2,981$) within the 176 STEC genomes sequenced.

- Long-read sequencing and bioinformatics processing produced 176 complete assemblies of STEC O157:H7. These genomes yielded 2,981 prophages in total of which 308 are *stx*-encoding prophages.
- Prophage content varied from 14 to 20 prophages per chromosome, with STEC O157:H7 belonging to lineage II having more prophages on average than those in lineage I and III (Figure 1).
- Prophages harbouring *stx1a* had a preferential *stx*-encoding bacteriophage integration site (SBI) of *yehV* (96.6%, 58/60) whereas *stx2c* had a preferential site of *sbcB* (98.25 112/114) and in most cases *stx2c* and *stx1a*-encoding showing little genetic diversity (Figure 2).
- *stx2a*-encoding prophages had the most diverse SBI sites, including *argW* (49.2% 66/134), *sbcB* (31.3% 42/134), *yecE* (8.2% 11/134) and this was reflected in terms of genetic distance as *stx2a*-encoding prophages showed much more diversity (Figure 2).
- There were 15 genomes that had multiple copies of *stx2a*-encoding prophages, inserted at different SBIs which would not have been detected via short-read sequencing.

DISCUSSION & CONCLUSIONS

Long-read sequencing allows for the characterisation of the accessory genome of STEC and thus allows us to detect and extract complete prophage sequences.

By cataloging the prophage content of STEC O157:H7 into a database is the first step into developing novel *in silico* methods of typing such as phage typing or determining if a given sample is domestic or non-domestic.

Accurately mapping the acquisition and loss of *stx*-encoding bacteriophages enables us to predict the virulence potential of the different STEC O157:H7 lineages, and to monitor for emerging threats to public health.

ACKNOWLEDGEMENTS

The research was funded by the National Institute for Health Research Protection Research Unit (NIHR HPRU) in Gastrointestinal Infections at University of Liverpool in partnership with UK Health Security Agency (UKHSA) formally Public Health England (PHE), in collaboration with University of Warwick. The views expressed are those of the author(s) and not necessarily the NIHR, the Department of Health and Social Care or UKHSA.

REFERENCES