



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Genome structural variation in Escherichia coli O157:H7

**Citation for published version:**

Fitzgerald, SF, Lupolova, N, Shaaban, S, Dallman, TJ, Greig, D, Allison, L, Tongue, SC, Evans, J, Henry, MK, McNeilly, TN, Bono, JL & Gally, DL 2021, 'Genome structural variation in Escherichia coli O157:H7', *Microbial Genomics*, vol. 7, no. 11. <https://doi.org/10.1099/mgen.0.000682>

**Digital Object Identifier (DOI):**

[10.1099/mgen.0.000682](https://doi.org/10.1099/mgen.0.000682)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**Published In:**

Microbial Genomics

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Genome structural variation in *Escherichia coli* O157:H7

Stephen F. Fitzgerald<sup>1,\*</sup>, Nadejda Lupolova<sup>1</sup>, Sharif Shaaban<sup>1</sup>, Timothy J. Dallman<sup>2</sup>, David Greig<sup>2</sup>, Lesley Allison<sup>3</sup>, Sue C. Tongue<sup>4</sup>, Judith Evans<sup>4</sup>, Madeleine K. Henry<sup>4</sup>, Tom N. McNeilly<sup>5</sup>, James L. Bono<sup>6</sup> and David L. Gally<sup>1,\*</sup>

## Abstract

The human zoonotic pathogen *Escherichia coli* O157:H7 is defined by its extensive prophage repertoire including those that encode Shiga toxin, the factor responsible for inducing life-threatening pathology in humans. As well as introducing genes that can contribute to the virulence of a strain, prophage can enable the generation of large-chromosomal rearrangements (LCRs) by homologous recombination. This work examines the types and frequencies of LCRs across the major lineages of the O157:H7 serotype. We demonstrate that LCRs are a major source of genomic variation across all lineages of *E. coli* O157:H7 and by using both optical mapping and Oxford Nanopore long-read sequencing prove that LCRs are generated in laboratory cultures started from a single colony and that these variants can be recovered from colonized cattle. LCRs are biased towards the terminus region of the genome and are bounded by specific prophages that share large regions of sequence homology associated with the recombinational activity. RNA transcriptional profiling and phenotyping of specific structural variants indicated that important virulence phenotypes such as Shiga-toxin production, type-3 secretion and motility can be affected by LCRs. In summary, *E. coli* O157:H7 has acquired multiple prophage regions over time that act to continually produce structural variants of the genome. These findings raise important questions about the significance of this prophage-mediated genome contingency to enhance adaptability between environments.

## DATA SUMMARY

All sequence files are available from the NCBI database.  
Accession numbers:

CP018237, CP015832, CP018239, CP018241, CP018243, CP018245, CP015831, CP018247 CP018252, CP018250, SAMN05544763, NC\_011353, CP008957, NC\_002695, CP010304, NC\_013008, CP038414, CP038402, CP038372, CP038366, CP038360, CP038357, CP038353, CP038349, CP038344, CP038416, CP038342, CP038309, CP039834, CP038333, CP038292, CP038290, CP038302, CP038300, CP038346, CP038355, CP038351, CP038282, CP038339, CP038319, CP038284, CP062749, CP062746, CP062744, CP062742, CP062736, CP062733, CP062731, CP062780, CP062729, CP062727, CP062725, CP062723, CP062721,

CP062719, CP062717, CP062715, CP062713, CP062711, CP062708, CP062705, CP062702, CP062700, CP062778, CP062774, CP062771, CP062769, CP062766, CP062761, CP062758, CP062755, CP062752, CP062763, CP062782, CP062739.

Short read archive files for strain sequenced in this study: PacBio - SRR15415552, SRR15415530, SRR15415529, SRR15415527, SRR15415523, SRR15415522, SRR15415519, SRR15415517, SRR15415604, SRR15415563, SRR15415603, SRR15415589, SRR15415598, SRR15415597, SRR15415605, SRR15415561, SRR15415588, SRR15415587, SRR15415606, SRR15415521, SRR15415518, SRR15415520, SRR15415600, SRR15415592, SRR15415558, SRR15376182, SRR15376181, SRR15376170, SRR15376159, SRR15376148, SRR15376140, SRR15376139,

Received 05 March 2021; Accepted 03 September 2021; Published 09 November 2021

**Author affiliations:** <sup>1</sup>Division of Infection and Immunity, The Roslin Institute and R(D)SVS, The University of Edinburgh, Easter Bush, Midlothian, EH25 9RG, UK; <sup>2</sup>Gastrointestinal Bacterial Reference Unit, 61 Colindale Avenue, Public Health England, NW9 5EQ London, UK; <sup>3</sup>Scottish *E. coli* O157/VTEC Reference Laboratory, Department of Laboratory Medicine, Royal Infirmary of Edinburgh, 51 Little France Crescent, Edinburgh EH16 4SA, UK; <sup>4</sup>Epidemiology Research Unit (Inverness), Department of Veterinary and Animal Science, Northern Faculty, Scotland's Rural College (SRUC), Scotland, IV2 5NA, UK; <sup>5</sup>Moredun Research Institute, Pentlands Science Park, Bush Loan, Penicuik, EH26 OPZ, UK; <sup>6</sup>United States Department of Agriculture, Agricultural Research Service, US Meat Animal Research Center, Clay Center, Nebraska, USA.

\*Correspondence: David L. Gally, dgally@ed.ac.uk; Stephen F. Fitzgerald, stephen.fitzgerald@roslin.ed.ac.uk

**Keywords:** O157:H7; cattle; duplication; *E. coli*; genome structure; inversion; optical mapping; PFGE; prophage; Shiga toxin; type 3 secretion.

**Abbreviations:** IS, insertion sequence; LB, Lysogeny Broth; LCR, large chromosomal rearrangement; ONT, Oxford Nanopore Technology; PT, phage type; Stx, Shiga toxin; SV, structural variant; Ter, terminus; T3SS, type-3 secretion system.

**Data statement:** All supporting data, code and protocols have been provided within the article or through supplementary data files. Two supplementary tables and eight supplementary figures are available with the online version of this article.

000682 © 2021 The Authors



This is an open-access article distributed under the terms of the Creative Commons Attribution License.

SRR15376138, SRR15376137, SRR15376136, SRR15376180, SRR15376179, SRR15376178, SRR15376177, SRR15376176, SRR15376175, SRR15376174, SRR15376173, SRR15376172, SRR15376171, SRR15376169, SRR15376168, SRR15376167, SRR15376166, SRR15376165, SRR15376164, SRR15376163, SRR15376162, SRR15376161, SRR15376160, SRR15376158, SRR15376157, SRR15376156, SRR15376155; Minion - SRR14434006; (Illumina) - SRR15376146, SRR15376144, SRR15376142, SRR15376141, SRR15376154, SRR15376153, SRR15376152, SRR15376149, SRR15376147, SRR15376145, SRR15376150, SRR15376151, SRR15376143. All defined in Table S1 (available in the online version of this article).

The RNAseq data is available in the Geo data base as study number; GSE158899. The remaining relevant data is within the manuscript and its supplementary data files.

## INTRODUCTION

Prophage are bacterial viruses integrated into the chromosome of their host and are major drivers of bacterial genome evolution, host and niche adaptation and virulence [1–3]. Prophage integration directly benefits the bacterial host by conferring resistance against some other lytic viruses [4], by carriage of virulence factors, including toxins and effector proteins [1, 5], enzymes involved in stress resistance [6] and the expression of both gene regulators and sRNAs capable of influencing the host gene regulatory network [2, 7]. Here we examine the impact prophages have on the structure of the bacterial genome through the generation of large-chromosomal rearrangements (LCRs).

*Escherichia coli* O157:H7 is a significant human zoonotic pathogen originating from ruminant hosts, especially cattle [8]. Over evolutionary time, numerous prophages (typically 16–25) have integrated into the genomes of *E. coli* O157:H7 strains with an integration bias towards the terminus (Ter) of replication [9]. Acquisition of these prophages, many of which are closely related  $\lambda$ -like phages, has driven the evolution of this pathogen by carriage of virulence genes including secreted effector proteins, sRNAs involved in virulence gene regulation and [7, 10], importantly, these prophages include those that encode Shiga-toxin (Stx) subtypes. Stx toxins are the main mediators of vascular endothelial cell killing in infected humans [11] and the resulting damage can lead to haemolytic uremic syndrome (HUS), often fatal, or lead to life-long kidney and brain damage [12–14]. *E. coli* O157:H7 strains divide into three phylogenetically distinct lineages (I, I/II and II). Those that represent a serious threat to human health belong to either lineage I or lineage I/II and the majority encode two sub-types of Stx:Stx2a and Stx2c. Stx2a is generally associated with more serious disease [11, 15–18] and the emergence of *E. coli* O157:H7 as a zoonotic threat correlates with the introduction of Stx2a-encoding prophage into the *E. coli* O157:H7 cattle population approximately 50 years ago [11].

There is published evidence that *E. coli* O157:H7 type strain EDL933 can undergo large-chromosomal rearrangements (LCRs), mainly inversions [19, 20], with these rearrangements being flanked by prophages. LCRs, such as inversions,

## Impact Statement

*Escherichia coli* has an 'open genome' and has acquired genetic information over evolutionary time, often in the form of bacteriophages that integrate into the bacterial genome (prophages). *E. coli* O157:H7 is a clonal serotype that is found primarily in ruminants such as cattle but can cause life-threatening infections in humans. *E. coli* O157:H7 isolates contain multiple prophages including those that encode Shiga-like toxins, which are responsible for the more serious disease associated with human infections. We show in this study that many of these prophages exhibit large regions of sequence similarity that allow rearrangements to occur in the genome generating structural variants. These occur routinely during bacterial culture in the laboratory and the variants are detected during animal colonization. The variants generated can give the bacteria-altered phenotypes, such as increased motility or toxin production and therefore represent a highly dynamic mechanism to generate variation in bacterial populations without a change in overall gene content.

uplications and translocations, occur by homologous recombination between repeat sequences on the same chromosome [21]. While LCRs arising between ribosomal *rrn* operons, pathogenicity islands and insertion sequence (IS) elements have been associated with speciation, diversification, outbreaks and immune evasion in bacteria [1, 22] few studies have examined LCRs arising from inter-prophage recombination and their impact on phenotype.

In this study we demonstrate that prophage-mediated LCRs are a major source of genomic variation across all lineages of *E. coli* O157:H7. We show that alternate chromosomal conformations are generated during laboratory culture and can be detected during host colonization. Specific LCRs were associated with changes in virulence phenotypes and we therefore propose that the generation of LCRs within *E. coli* O157:H7 populations can lead to phenotypic heterogeneity with the potential for selection of specific variant sub-populations in the different niches encountered by the bacteria.

## METHODS

### Bacterial strains and culture conditions

Bacterial strains and plasmids used in this study are listed in Table S1. Bacteria were cultured in Lysogeny Broth (LB) or M9 minimal media (Sigma-Aldrich) supplemented with 0.2% glucose, 2 mM MgSO<sub>4</sub> and 0.1 mM CaCl<sub>2</sub>. For TTSS expression bacteria were cultured overnight in LB and then inoculated into minimal essential medium (MEM)-HEPES (Sigma-Aldrich) supplemented with 0.1% glucose and 250 nM Fe(NO<sub>3</sub>)<sub>3</sub>. Antibiotics were used at the following concentrations when required: Chloramphenicol (50  $\mu$ g ml<sup>-1</sup>), Mitomycin C (2  $\mu$ g ml<sup>-1</sup>), Nalidixic acid (50  $\mu$ g ml<sup>-1</sup>).

## PacBio long-read sequencing

A total of 72 whole-genome sequences were used for analysis in this study. Long-read sequencing of 49 strains was carried out for this study and the remaining were publicly available in the National Centre for Biotechnology Information (NCBI) database (Table S1).

Sequencing of the isolates was conducted using a PacBio RS II long-read sequencing platform and carried out at the U.S. Department of Agriculture sequencing core facility in Clay Center, Nebraska, USA. Qiagen Genomic-tip 100 G<sup>-1</sup> columns and a modified protocol, as previously described [23], were used to extract high molecular weight DNA. Using a g-TUBE (Corvaris), 10 µg of DNA was sheared to a targeted size of 20 kb and concentrated using 0.45× volume of AMPure PB magnetic beads (Pacific Biosciences). Following the manufacturer's protocol, 5 µg sheared DNA and the PacBio DNA SMRTbell Template Prep kit 1.0 were used to create the sequencing libraries. A BluePippin instrument (Sage Science) with the SMRTbell 15–20 kb setting was used to size select 10 kb or larger fragments. The library was bound with polymerase P5 and sequencing was conducted with the C3 chemistry and the 120 min data collection protocol. Individual libraries were constructed from some of the strain DNA preparations described above using an Illumina Nextera XT DNA sample preparation kits with appropriate indices tags according to the manufacturer's instructions (Illumina, San Diego, CA). The libraries were pooled together and run on an Illumina MiSeq DNA sequencer (Illumina, San Diego, CA). The genome of each strain was sequenced to a targeted depth of 50× coverage.

## Genome assembly and annotation

SMRT analysis was used to generate a FASTQ file from the PacBio reads, which were then error-corrected using PbcR with self-correction [24]. The Celera Assembler was used to assemble the longest 20× coverage of the corrected reads. The resulting contigs were improved using Quiver [25] and annotation was conducted using a local instance of Do-It-Yourself Annotator (DIYA) [26]. Geneious (Biomatters) was used to remove duplicated sequence from the 5' and 3' ends to generate the circularized chromosome. To correct PacBio sequencing errors (homopolymers and SNPs), Illumina reads were mapped to the Quiver-polished chromosome using Pilon [27]. Then, both PacBio and Illumina reads were mapped to the Pilon-generated chromosome using Geneious Mapper. Additional sequencing errors were identified and corrected by manual editing in Geneious, resulting in a finished closed circularized chromosome. OriFinder was used to determine the origin of replication [28] and the chromosome was reoriented using the origin as base number one. Prophage regions were identified as described previously [9] using PHASTER [29].

## MinION sequencing and SV read detection

Strain 9000 was sequenced by Oxford nanopore technologies MinION sequencing. High molecular weight genomic DNA was extracted from strain 9000 grown in LB (OD<sub>600</sub>=0.7) by standard phenol:chloroform extraction [30]. Genomic DNA was

purified using Qiagen G100 Genomic Tips (Qiagen) with minor alterations including no vigorous mixing steps and final elution in 100 µl of nuclease free water and quantified using a Qubit and the HS (high sensitivity) dsDNA assay kit (ThermoFisher Scientific), following the manufacturer's instructions. Library preparation was performed using the Ligation kit SQK-LSK109 (Oxford Nanopore Technologies). The prepared libraries were loaded onto a FLO-MIN106 R9.4.1D flow cell (Oxford Nanopore Technologies) and sequenced using the MinION (Oxford Nanopore Technologies) for 72 h. Data produced in a raw FAST5 format was basecalled and de-multiplexed using Guppy v3.2.4 using the FAST protocol (Oxford Nanopore Technologies) into FASTQ format.

To determine if the Nanopore-sequenced strain 9000 contained reads supporting multiple isoforms of the chromosome, Minimap2 v2.17 [31] and Samtools v1.7 [32] were used to align the Nanopore reads (removing secondary aligning reads) to samples Z1615, Z1723 and Z1767 each representing a different chromosomal isoform. Using Samtools v1.7 [32] and Bedtools v2.29.2 [33] reads were identified at either end of the each of the 5' and 3' breakpoints identified in those conformations. The number of reads that crossed each end of the 5' and 3' breakpoints for both conformations was calculated again using Samtools v1.7 [32] and Bedtools v2.29.2 [33].

## Whole-genome comparisons

Pairwise whole-genome alignments were conducted with Easyfig [34] as described previously [9]. Genome .gbk files were modified so that prophages were represented as coloured blocks. AvrII restriction sites were identified in selected genomes using UGENE [35] and their loci were added to the respective genome .gbk files. Pairwise whole-genome alignments reference genomes from each lineage (9000, Sakai, TW14359 and 180) and each genome were performed using BLASTn [36] with the following parameters (-evalue 1e-10 -best\_hit\_score\_edge 0.05 -best\_hit\_overhang 0.25 -perc\_identity 70 -max\_target\_seqs 1 -outfmt 6). From the resulting alignment files, LCRs were identified within each genome by filtering all inverted homologous regions ≥50000 bp relative to each reference strain.

## Mapping homologous regions

Homologous regions within each genome were identified using BLASTn [36]. BLASTn was performed on each individual genome using the same genome sequence as both reference and query with the following parameters (-evalue 1e-10 -best\_hit\_score\_edge 0.05 -best\_hit\_overhang 0.25 -perc\_identity 98 -max\_target\_seqs 1 -outfmt 6). Homologous regions that satisfied three conditions simultaneously were extracted from the BLAST output: (1) homologous regions were ≥5000 bp; (2) homologous regions ≥5000 bp were present in the genome at a frequency ≥2; (3) homologous regions were located before and after *dif* (terminus of replication). Equivalent analysis was repeated to determine homologous regions ≥8000 bp. Significant differences in the total number of repeats detected



between lineages and the bias of repeat regions toward Ter was determined by one-way ANOVA with Dunnett's multiple comparisons test.

Circos plots [37] were used to visualize linked regions of homologous sequence within the genomes of selected strains. Custom circos input files were generated in which the data matrix was modified such that each circular genome was divided at prophage boundaries. Linked homologous regions and their sizes were determined using BLAST scores derived when querying a selected genome sequence to itself. Only BLAST hits with  $\geq 98\%$  sequence homology and that were  $\geq 5000$  bp in length were included in the data matrix of circos input files. Within Circos plots the width of linked segments is proportional to the length of BLAST hits. Circos does not exactly map homology hits to linked chromosomal/prophage regions, instead connecting segments originate and end at the earliest available location within the linked region.

### Phylogeny of the strains

A core gene alignment was extracted from the fully assembled and annotated PacBio genomes of 72 strains using ROARY [38] with parameters (-e -n -r -s -ap). The extracted multiple alignment was used to generate maximum-likelihood phylogenetic trees using FastTree [39] (-gtr) and trees were visualized with iTOL [40]. To determine the phylogenetic relationship of PT21/28 strains high-quality illumina sequencing reads were mapped to the reference STEC O157:H7 strain, Sakai (GenBank accession BA000007), using Burrows-Wheeler Aligner – Maximum Exact Matching [BWA MEM (v0.7.2)] [41]. The sequence alignment map output from BWA were sorted and indexed to produce a binary alignment map (BAM) using Samtools (v1.1) [42]. Genome Analysis Toolkit (GATK v2.6.5) was then used to create a variant call format (VCF) file from each of the sorted BAMs, which were further parsed to extract only SNP positions of high quality [mapping quality (MQ) >30, depth (DP) >10, variant ratio >0.9]. Hierarchical single linkage clustering was performed on the pairwise SNP difference between all isolates at descending distance thresholds ( $\Delta 250$ ,  $\Delta 100$ ,  $\Delta 50$ ,  $\Delta 25$ ,  $\Delta 10$ ,  $\Delta 5$ ,  $\Delta 0$ ) [43]. SNP alignments were created tolerating positions where >80% of isolates had a base call with regions of recombination masked using Gubbins v2.0.0 [44]. Maximum-likelihood phylogenies were computed using IQ-TREE v2.0.4 [45] with the best-fit model automatically selected and near zero branches collapsed into polytomies.

### Pulsed-field gel electrophoresis

All strains analysed by PFGE were cultured in LB or M9 medium at 37 °C overnight with agitation. Genomic DNA was purified using the CHEF Bacterial Genomic DNA Plug Kit (Bio-Rad) according to manufacturer guidelines. DNA restriction digestion with AvrII (BlnI) (Takara) and subsequent PFGE was done according to the PulseNet O157:H7 guidelines [46], using a CHEF-DR III system. *In silico* AvrII (BlnI) restriction digests of selected genomes were carried out in CLC Genomics Workbench (Qiagen).

### RNA sequencing

Total RNA was extracted from three biological replicates of strains 9000, Z1615 Z1723, Z1767 using mirVana miRNA Isolation Kit (ThermoFisher) according to manufacturer guidelines.

Strains were cultured in either LB or M9 media to  $OD_{600} = 0.7$ . Ribosome depletion, cDNA library preparation and Illumina sequencing was carried out by Vertis Biotechnologie AG (Freising, Germany). Total RNA samples were purified and concentrated using the Agencourt RNAClean XP kit (Beckman Coulter Genomics) and the RNA integrity was assessed by capillary electrophoresis. Ribosomal RNA molecules were depleted using the Ribo-Zero rRNA Removal Kit for bacteria (Illumina). The ribodepleted RNA samples were first fragmented using ultrasound (four pulses of 30 s each at 4 °C) and oligonucleotide adapters were then ligated to the 3' end of the RNA molecules. First-strand cDNA synthesis was performed using M-MLV reverse transcriptase and the 3' adapter as primer. The first-strand cDNA was purified and the 5' Illumina TruSeq sequencing adapter was ligated to the 3' end of the antisense cDNA. The resulting cDNA was PCR-amplified to about 10–20 ng  $\mu\text{l}^{-1}$  using a high-fidelity DNA polymerase. The cDNA was purified using the Agencourt AMPure XP kit (Beckman Coulter Genomics) and was analysed by capillary electrophoresis. Purified cDNA was pooled and sequenced on an Illumina NextSeq 500 system using 75 bp read length. RNA-sequencing reads were mapped to the strain 9000 reference genome (CP018252.1) using STAR 2.7.0e [47] with the following parameters (--quantMode GeneCounts and --sjdbGTFfeatureExon CDS). Prior to read mapping the reference strain 9000 was annotated using Prodigal version 2.6 [48]. The loci of previously identified *E. coli* O157:H7 sRNAs [7] were found in strain 9000 using BLASTn and manually added to strain 9000 .gtf file. Column 3 of the reference GTF file (feature) was manually modified to CDS for all genetic features. Differential expressed (DE) genes were identified with edgeR [49] ( $P$ -values=0.05) using the glmQLFit +glmQLFTest parameters. RNA-seq data was uploaded to NCBI Gene Expression Omnibus (GEO) (Accession: **GSE158899**).

### Stx toxin ELISA

In total, 3 ml LB was inoculated directly from glycerol stocks and grown overnight at 37 °C. 6 ml LB was inoculated 1/100 from overnight cultures and grown to an  $OD_{600\text{ nm}} = 0.6$ –0.8. Mitomycin C ( $2\ \mu\text{g ml}^{-1}$ ) was added and lysis allowed to proceed for 24 h. After 24 h, 1 ml culture was taken and live cells and cell debris removed by centrifugation (13000 r.p.m.). Stx toxin containing supernatants were further sterilized by syringe filtering ( $0.22\ \mu\text{m}$ ; Milipore). The level of Stx toxin in each sample was assayed using the RIDASCREEN Verotoxin ELISA kit (R-Biopharm) according to manufacturer guidelines. Differences Stx2 production was assessed by ordinary one-way ANOVA with multiple comparisons where each strain was compared with Z1723.

### Stx Vero cell toxicity

Cytotoxicity of Stx2 toxin was measured on Vero cell monolayers cultured in RPMI medium (Sigma-Aldrich). Cells (100  $\mu$ l) were plated into 96-well microtitre plates and at ~75% confluence the culture medium was replaced with RPMI medium containing diluted (1 : 1000) Stx2 toxin supernatants. Vero cells were exposed to Stx2 toxin for 72 h at 37 °C, 5% CO<sub>2</sub>. Surviving cells were fixed using paraformaldehyde (2%) and stained with crystal violet (10%). Crystal violet was solubilized with 10% acetic acid live/dead cells were quantified spectrophotometrically at 590 nm. Cells exposed to Triton X-100 (0.1%) and RPMI were used as positive and negative controls for toxicity, respectively. Strain toxicity was expressed as a percentage of the toxicity measured for RPMI control. Strain toxicity was analysed by ordinary one-way ANOVA with multiple comparisons where each strain was compared with Z1723.

### Fitness assays

The fitness of strain 9000 variants grown in M9 media was calculated as described previously [50, 51]. Viable-cell counts for each competing strain were determined at time zero ( $t=0$ ) and again after 24 h of co-culturing by selective plating. Fitness was calculated using the formula:

$$\text{Fitness index (f.i.)} = \frac{\text{LN} [N_i(1) / N_i(0)]}{\text{LN} [N_j(1) / N_j(0)]},$$

where  $N_i(0)$  and  $N_i(1)$ =initial and final colony counts of strain Z1723 or 9000, respectively, and

$N_j(0)$  and  $N_j(1)$ =initial and final colony counts of structural variant strain (Z1766, Z1767 or Z1615), respectively.

For controls WT strain 9000 or Z1723 were competed against NaI<sup>r</sup> derivatives generated previously [52]. Fitness was analysed by ordinary one-way ANOVA with Dunnett's multiple comparisons test where each strain was compared with the control.

### RT-qPCR

Total RNA was extracted from cell pellets using a RNeasy Mini kit (Qiagen) according to manufacturer guidelines. Extracted RNA was quantified and 2  $\mu$ g of each sample was DNase treated using TURBO DNA-free kit. 200 ng of DNase-treated RNA was then converted to cDNA using iScript Reverse Transcription Supermix (Bio-Rad) according to manufacturer guidelines. All qPCR reactions were carried out using iQ Syber Green supermix (Bio-Rad) and *stx2a* specific primers (IDT-DNA): stx2a-F-GAAGAAGATGTTTATG GCGGTTT, stx2a-R-CCCGTCAACCTTCACTGTAA. Cycling conditions were 95 °C for 15 s (1 cycle), 95 °C for 15 s; 60 °C for 1 min (40 cycles). Gene expression was quantified relative to a standard curve generated from Z1723 genomic DNA.

### Optical mapping

Strains Z1723 and Z1767 were cultured from a single colony in LB medium to an OD<sub>600</sub>=0.7. 1 ml of cells/agarose plug

were harvested (4000 g, 5 min) and intact chromosomes were extracted according to the Bionano Prep Cell Culture DNA Isolation Protocol (Bionano). Briefly, harvested cells were washed twice in Bionano Cell Buffer (Bionano). Washed cells were embedded in 2% low-melt agarose plugs and cells were lysed (1 h at 37 °C) with lysozyme enzyme (100  $\mu$ l) using CHEF Bacterial Genomic DNA Plug Kit (Bio-Rad). DNA containing plugs were washed twice with nuclease free water then treated with Proteinase K (Qiagen) in Bionano Lysis Buffer according to the Bionano Prep Cell Culture DNA Isolation Protocol. All subsequent procedure steps (RNase treatment, DNA extraction, quantitation and labelling) and optical mapping on Bionano Irys platform were provided as a service by Earlham Institute (Norwich, UK). Structural variant analysis was provided by Bionano and structural variant maps visualized using Bionano access (Bionano).

### TTSS expression and secretion

Expression of *ler* and *sepL* was measured using GFP reporter fusion plasmids pDW-LEE1 [53] and pDW6 [54], respectively. Reporter plasmids were transformed into strains Z1723, Z1766 and Z1767 by electroporation and transformants were cultured overnight in LB media supplemented with chloramphenicol (50  $\mu$ g ml<sup>-1</sup>). Overnight cultures were diluted 1 : 100 into MEM-HEPES and grown at 37 °C (200 r.p.m.) to an OD<sub>600</sub> 0.8–1.0. GFP fluorescence of 200  $\mu$ l aliquots was measured in a 96-well blank microtitre plate using a FLUOstar Optima plate reader (BMG, Germany). The Gfp promoter-less plasmid pKC26 was used as a control [55].

For EspD secretion, bacteria were cultured in 50 ml of MEM-HEPES at 37 °C (200 r.p.m.) to an OD<sub>600</sub> of 0.8–1.0. Bacterial cells were pelleted by centrifugation at 4000 g for 20 min, and supernatants were passed through low protein binding filters (0.45  $\mu$ m). Then, 10% TCA was used to precipitate proteins overnight, which were separated by centrifugation at 4000 g for 30 min at 4 °C. The proteins were suspended in 150  $\mu$ l of 1.5 M Tris (pH 8.8). For bacterial lysates, bacterial pellets were suspended directly in SDS PAGE loading buffer. Proteins were separated by SDS-PAGE using standard methods and Western blotting performed as described previously for EspD and RecA [56].

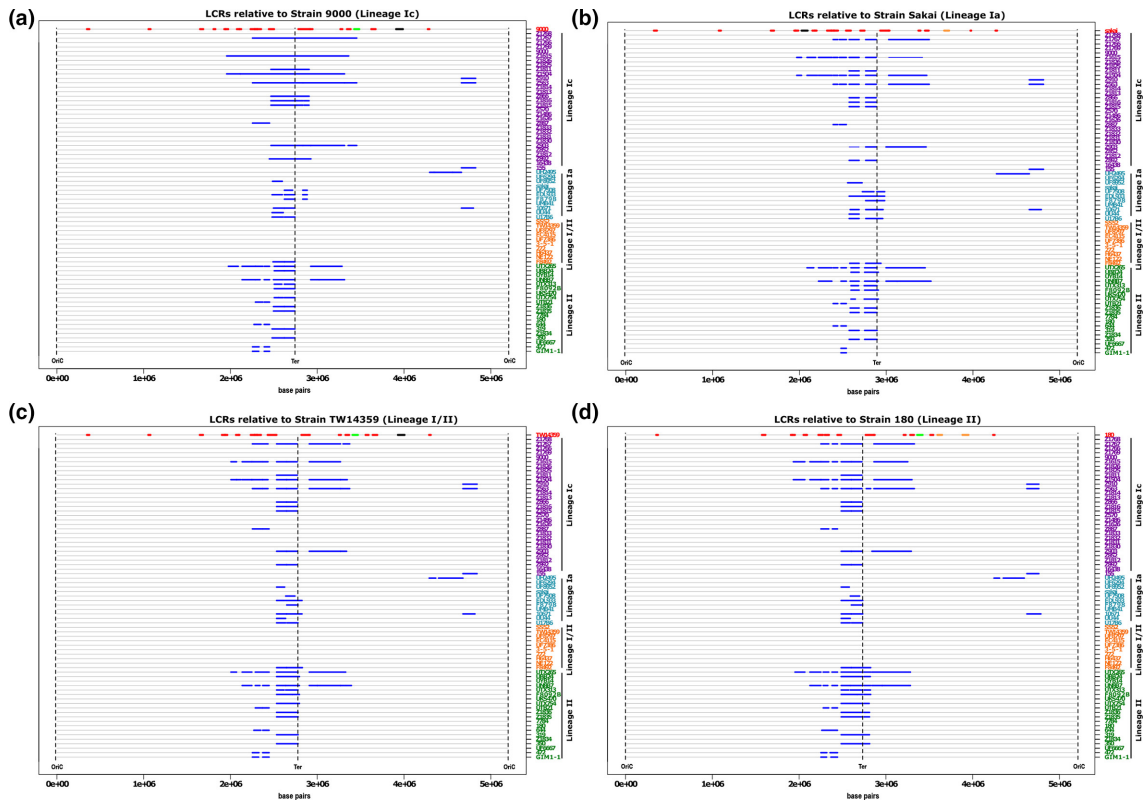
### Motility assay

The motility of strains was determined using Tryptone swarm plates as described previously [57]. The centre of Tryptone swarm plates were inoculated with 2  $\mu$ l volumes (OD<sub>600</sub>=1.0) of each strain being assessed and incubated for 6 h at 37 °C. The swarm radius after 6 h was measured manually.

## RESULTS

### LCRs shape the *E. coli* O157:H7 genome

To examine the extent of genomic diversity generated by LCRs in the *E. coli* O157:H7 clonal group, we examined the whole-genome sequences of 72 isolates, the majority of which were generated by PacBio long-read sequencing (Table



**Fig. 1.** Position of major chromosomal rearrangements in *E. coli* O157:H7 genomes. The relative positions of all LCRs  $\geq 50$  kb (blue lines) are marked on the chromosomal maps (grey line) of strains from lineage Ic (purple), lineage Ia (blue), lineage I/II (orange) and lineage II (green). Strains are sorted based on their phylogenetic relatedness. Chromosomes are centred by the replication terminus (Ter), beginning and ending at the origin of replication (OriC). LCRs are shown relative to four reference strains: (a) 9000, lineage Ic; (b) Sakai, lineage Ia; (c) TW14359, lineage I/II; (d) 180, lineage II. The loci of non-*stx* prophage (red line) and prophage carrying *stx1* (orange), *stx2c* (green) and *stx2a* (black) are mapped for each reference strain.

S1). Strains analysed were representative of the main *E. coli* O157:H7 lineages (I, I/II and II) and included multiple sub-Lineage Ic, phage type (PT) 21/28 isolates, which have been responsible for the majority of serious human infections in the UK over the last two decades [11]. This genome dataset included previously sequenced complete genomes from each lineage, including strains Sakai (NC\_002695.2), EDL933 (CP008957.1) and TW14359 (CP001368) (Table S1).

Pairwise alignment of all 72 genomes identified LCRs, predominantly large inversions, as a common source of genomic variation between isolates within each *E. coli* O157:H7 lineage with the exception of lineage I/II (Fig. S1). In addition, each genome was individually aligned against a representative reference strain from each of four lineages and the chromosomal loci of all LCRs  $> 50$  kb were mapped (Fig. 1a–d). The reference strains were Strain 9000 (lineage I c), Sakai (lineage I a), TW14359 (lineage I/II) and Strain 180 (lineage II).

LCRs  $> 50$  kb were frequently identified in lineages Ia, Ic and II irrespective of the reference strain used for alignment, however it was evident that lineage I/II strains exhibited less variation (Fig. 1). Strains from lineage Ic and lineage II

exhibited the most variation at this macro level with an average of 43 and 37 LCRs identified, respectively (Table 1, Fig. 1). Strains from lineage Ia were less variable with an average of 14 LCRs identified across all strains and the least genomic variation with respect to the reference strains was observed for strains from lineage I/II with an average of just 2.5 LCRs identified in a single strain, F8492. We note that strain F8492 was a singleton isolate that grouped closely with our other representative lineage I/II strains (Fig. S2). Lineage I/II strains were also the least variable when the number of LCRs identified were corrected to account for the unequal number of strains analysed within each lineage (Table 1). To examine this further, we plotted the average size of all LCRs with a lower cut-off of  $> 20$  kb that could be detected in each strain relative to the four reference genomes (Fig. S2a–d). At this lower cut-off, LCRs ranging between 20 kb and 30 kb were identified in lineage I/II that were generally consistent across all lineage I/II strains relative to each reference genome. These results indicate that the macro genome conformation of lineage I/II strains is highly conserved. While LCRs can occur within lineage I/II strains they have a reduced capacity to generate larger LCRs  $> 30$  kb compared with the two other lineages.

**Table 1.** Mean number and size of LCRs relative to each reference genome

Reference	Total number of LCRs			
	Lineage Ia	Lineage Ic	Lineage I/II	Lineage II
9000	18	36	2	43
Sakai	15	47	2	38
TW14359	13	47	2	34
180	13	43	4	41
Mean	14.75	43.25	2.5	39
<b>LCRs per lineage/strain</b>				
9000	1.64	1.16	0.2	2.15
Sakai	1.36	1.52	0.2	1.9
TW14359	1.18	1.52	0.2	1.7
180	1.18	1.39	0.4	2.05
Mean	1.34	1.4	0.25	1.95
<b>Average LCR size (bp)</b>				
9000	109705	376229	126954	110504
Sakai	151151	143497	142760	114507
TW14359	155589	134568	149038	127925
180	131237	144404	119660	128622
Mean	136920	199675	134603	120390

For all lineages, LCRs were biased toward the chromosomal terminus of replication (Ter) with the majority located between 2–3.5 Mbp (Fig. 1). The largest LCR identified was a 1.4 Mbp inversion, which was detected in lineage Ic strain Z1615 (Figs 1 and S1). The average length of LCRs detected ranged between 109–376 Kbp depending on which reference strain was used for alignment with the largest LCRs detected within lineage Ic strains when aligned against lineage Ic reference strain 9000 (Table 1, Fig. S3a). Mapping the chromosomal position of prophages within each reference genome further demonstrated that most LCRs were bounded by prophages (marked in red in comparison strain, Fig. 1). Furthermore, many of the LCRs identified had prophage Stx2c ( $\Phi$ Stx2c) as a boundary, particularly those occurring within lineage Ic strains.

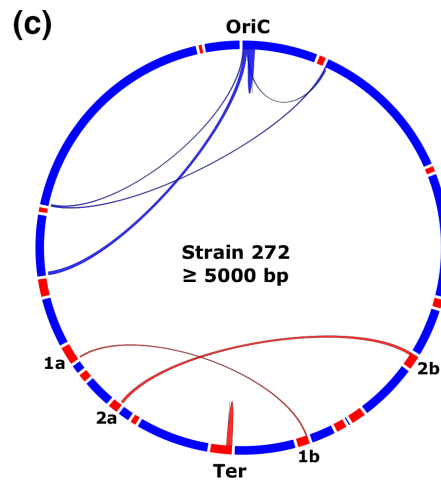
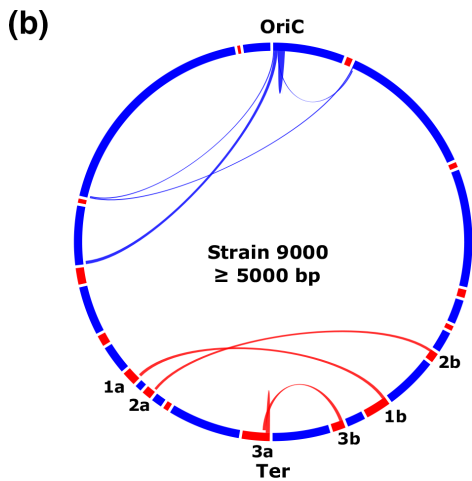
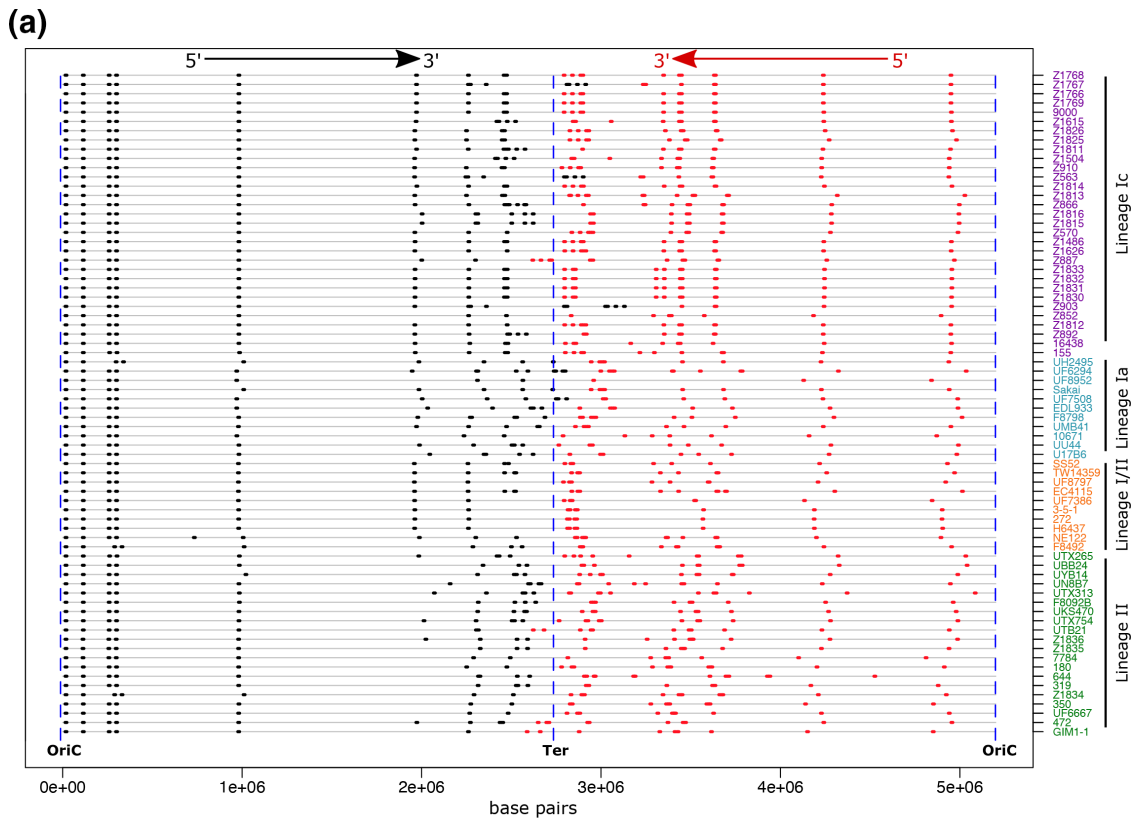
### LCRs map to repeated regions of homology in prophage

Mechanistically, chromosomal inversions typically involve recombination between inverted repeat regions of homologous sequences [22, 58]. As inversions were the dominant LCR identified in our analysis (Fig. 1), we mapped the chromosomal position and direction of all homologous regions for each *E. coli* O157:H7 strain (Fig. 2, Table S2). To avoid detection of the numerous insertion sequence (IS) elements present in *E. coli* O157:H7 genomes [59] we restricted our analysis to regions that shared  $\geq 98\%$  sequence homology, were  $\geq 5000$  bp

and occurred in the chromosome with a frequency  $\geq 2$ . Repeat regions were unequally distributed throughout the chromosome with a bias toward Ter and were conserved as inverted repeats at either side of Ter (Fig. 2a). When each genome was subdivided into 1 Mbp domains, significantly more repeats were located within the 2–3 Mbp domain ( $P < 0.0001$ ) adjacent to Ter than any other domain of the chromosome (Fig. S3b). Significantly more repeats were also located within the 3–4 Mbp domain ( $P < 0.0001$ ) adjacent to Ter than the 1–2 Mbp, 4–5 Mbp and 5–6 Mbp regions but not the 0–1 Mbp domain ( $P = 0.71$ ). All repeat regions identified in the terminal half of the chromosome mapped within prophage (Fig. 2b, c) and specific combinations of these repeated regions matched the boundaries for identified LCRs. For example, specific recombination between regions 1a and 1b of Strain 9000 in Fig. 2b would generate the LCR present in isogenic strain Z1767 and recombination between 2a and 2b would generate the LCR present in isogenic strain Z1615. These results indicate that homologous prophage sequences are hotspots for recombination resulting in the generation LCRs in *E. coli* O157:H7 strains.

It was evident that specific combinations of inverted repeat regions were present in the different lineages and sub-lineages of *E. coli* O157:H7 (Fig. 2a). We reasoned that the frequency of recombinational events would be greater in strains with more homologous repeat regions and *vice versa*. Indeed, strains





**Fig. 2.** Mapping homologous regions ( $\geq 5000$  bp) in *E. coli* O157:H7. The loci of all regions of homology  $\geq 5000$  bp (black/red) that are present as  $\geq 2$  copies per genome are mapped on the chromosomes (grey line) of strains from lineage Ic (purple), lineage Ia (blue), lineage I/II (orange) and lineage II (green) (a). The directions 5' – 3' of homologous sequences relative to OriC are shown with black indicating the inverse direction to red. Circos plots for lineage Ic strain 9000 (b) and lineage I/II strain 272 (c) show paired regions of homology. Prophage loci (red blocks) are shown on the respective circular genome maps (blue). Paired homologous regions are joined by arches: chromosomal (blue) and within prophage (red). Homologous prophage regions  $>5000$  bp are paired: 1a with 1b, 2a with 2b and 3a with 3b.

from lineage Ic, in which the greatest number of LCRs were identified (Table 1), had significantly more repeat regions  $\geq 5000$  bp ( $P < 0.05$ ) (Fig. S4a) and  $\geq 8000$  bp ( $P < 0.0001$ ) (Fig. S4b) than those from any other lineage. Conversely, lineage I/II strains, in which only a single LCR was identified, had

fewer homologous repeat regions  $\geq 5000$  bp than strains from any other lineage (Fig. S4a) and significantly less homologous repeat regions  $\geq 8000$  bp ( $P < 0.01$ ) (Fig. S4b).

## LCRs underpin PFGE type expansion in lineage Ic PT21/28 strains

In the UK, PT21/28 strains from lineage Ic have arisen as the dominant PT associated with severe human infections over the last 20 years [11]. Based on standard PFGE-typing methods, PT21/28 isolates have undergone a PFGE type expansion from an initial five PFGE types (profiles A - E) present in the UK in 1994 to >30 distinct PFGE profiles by 2013 (personal communication from Dr Lesley Allison Scottish *E. coli* reference laboratory-SERL) (Figs 3a and S5). LCRs were shown to generate changes in the PFGE type of strain EDL933 [19], we therefore determined if LCRs could be responsible for the expansion in PFGE patterns for PT21/28 strains described above. We sequenced ten PT21/28 isolates by PacBio long-read sequencing that differed in PFGE type. Strains were selected from throughout the PT21/28 core SNP-based phylogeny (Fig. S5) and the dataset included two isolates with identical SNP addresses (Z910 and Z563; zero SNP differences in the core genome) but with distinct PFGE profiles.

Sequence analysis showed that all ten strains differed by <70 SNPs in their core genomes (Fig. S5). Although examples of phage gain/loss ( $n=2$ ) were apparent, pairwise whole-genome comparisons showed that LCRs were the dominant source of genomic variation at the macro scale (Fig. 3b). Reference laboratories specializing in STEC diagnostics in the UK used *AvrII* and/or *XbaI* restriction enzymes when determining the PFGE type for an isolate. When all *AvrII* restriction sites were mapped in each isolate (Fig. 3b) it was evident that the loci of most sites were strongly conserved. However significant strain variation in *AvrII* loci was observed within the *Ter* region of the chromosome that was associated with LCRs. For example, strains Z910 and Z563, which were identical at the core SNP level, differed by a single 1.2 Mbp chromosomal inversion that involved recombination with  $\Phi$ Stx2c and resulted in the repositioning of four *AvrII* sites. Additional sequences containing *AvrII* sites were present in strains Z869 ( $n=2$ ) and Z570 ( $n=1$ ) that were not identified in other strains (Fig. 3b). The majority of *AvrII* loci variation and therefore PFGE type variation was generated by LCRs.

To confirm that the variation in *AvrII* loci generated by LCRs observed in our PacBio assemblies matched the actual chromosome configuration of each isolate we determined the PFGE profile for each strain after *AvrII* restriction digestion (Fig. 3c) and compared it to *in silico* *AvrII* digests of their respective PacBio assemblies (Fig. 3d). Both *in vivo* and *in silico* *AvrII* digestion patterns were matched for nine/ten strains analysed confirming the presence of those LCRs identified by PacBio long-read sequencing and the rearrangement of *AvrII* loci by these LCRs to generate different PFGE types. The exception was strain Z892 in which an unexplained digestion product was present after *in vivo* digestion that was not predicted from the PacBio sequence.

Based on these results we propose that the majority of the PFGE variation amongst PT21/28 strains, as depicted in Figs 3a and S5, is generated by LCRs. It was also evident

from the PFGE analyses that the strains cultured under these laboratory conditions had the majority of their genomes in a single confirmation (Fig. 3c). Of note, the most frequently occurring PT21/28 strain PFGE profile was type 'C' later defined as profile A\_11b (Figs 3a and S5). Phylogenetically, this specific profile re-occurs throughout the sub-lineage indicating that it is likely an ancestral confirmation or strains can repeatedly return to this chromosome conformation.

## Detection of LCRs during cattle colonisation

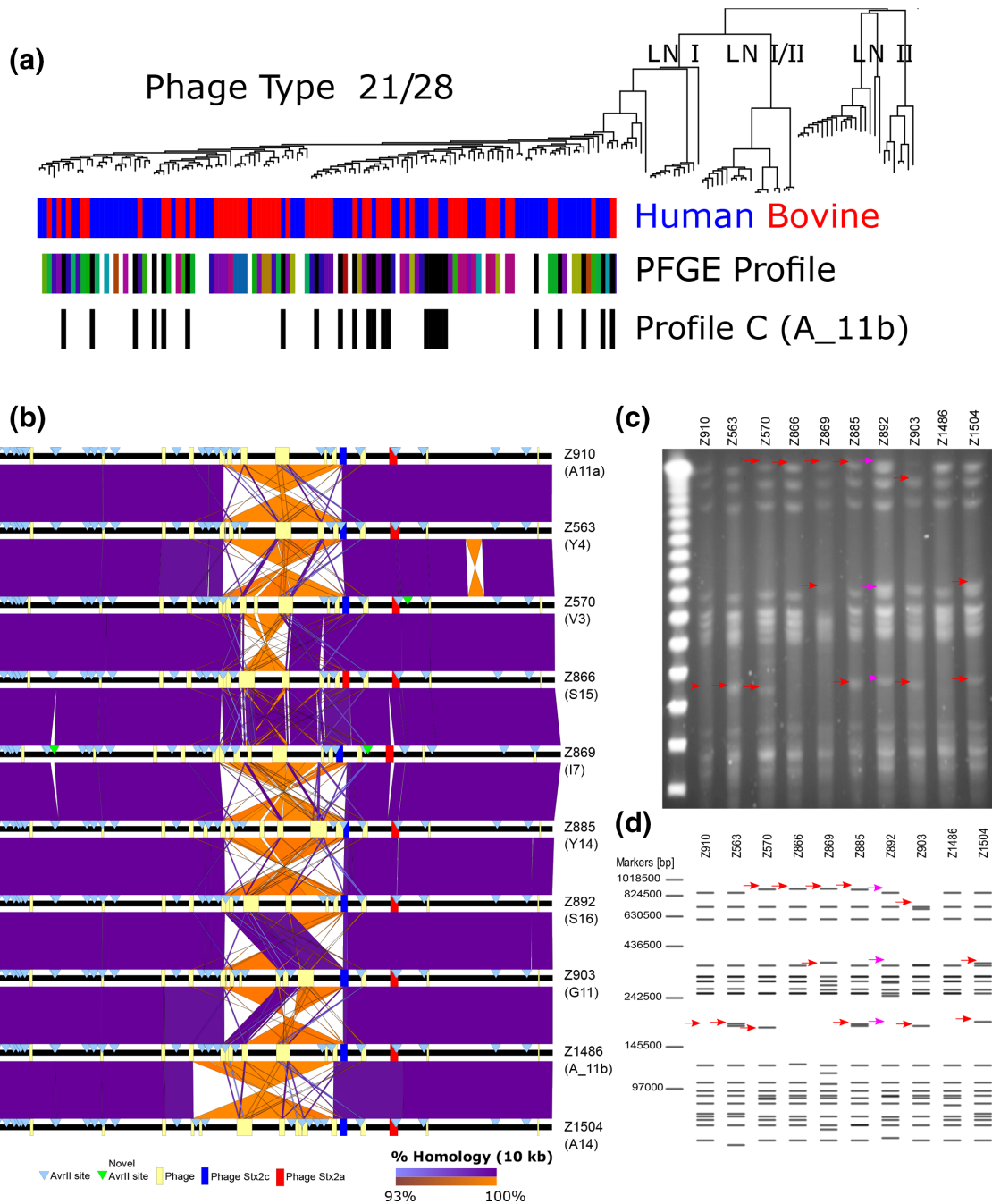
Previously, we carried out a series of published and unpublished *in vivo* cattle colonization studies focused on *E. coli* O157:H7 strain 9000 [52, 60]. We wanted to determine if LCR variants with different chromosome conformations to that of the inoculum, strain 9000, could be recovered from colonised cattle. For this we used both PacBio long-read sequencing and *AvrII* PFGE profiling of recovered colonies (Table S1). We note that recovery of variants from cattle does not strictly mean they were selected *in vivo* but does indicate that these variants can be recovered from the animal host.

Pairwise whole-genome comparisons of strains 9000 and Z1615 from trial 1 (Fig. 4a) showed a 1.4 Mbp inversion had occurred in derivative strain Z1615<sup>inv1.4</sup> relative to strain 9000. As outlined in Fig. 2(b) the boundaries of this LCR mapped to large inverted repeat sequences within prophage located either side of *Ter*. Distinct PFGE profiles were observed for strains 9000 and Z1615<sup>inv1.4</sup> following *AvrII* digestion, each matched their respective *in silico* *AvrII* digestion profiles (Fig. S6a) and confirmed the presence of the LCR identified in Z1615<sup>inv1.4</sup>. We analysed a further 11 recovered isolates from two experimental trials (trial 1 and trial 3), in which strain 9000 was the inoculum, by PFGE (Fig. S6b). Isolates were collected from a number of different animals and dates (Table S1). Three additional isolates of the 11 tested matched the PFGE profile of Z1615 containing the large 1.4 Mbp inversion.

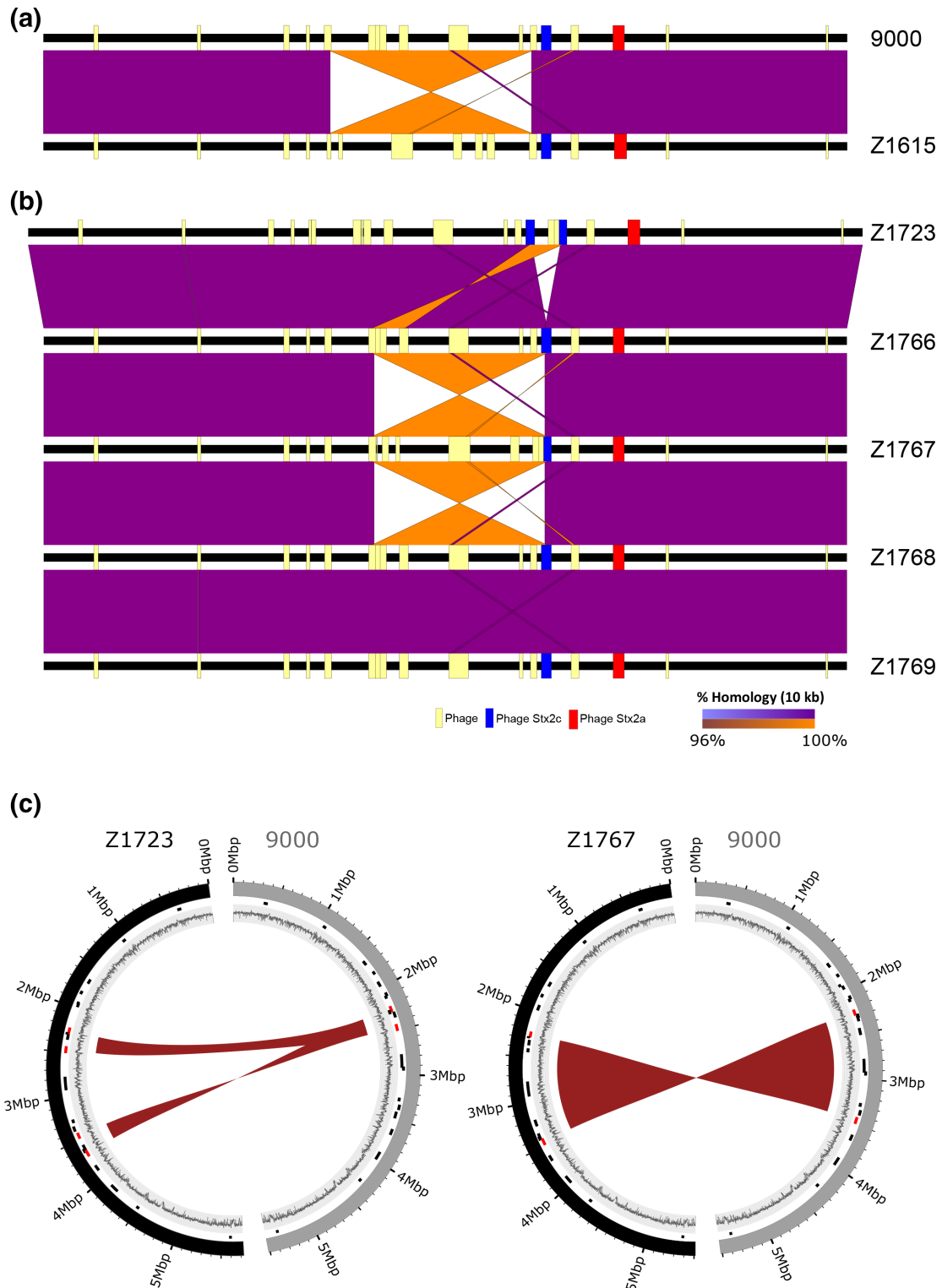
Two additional LCRs were identified from the five isolates examined from trial 2 (Fig. 4b) both of which involved recombination with the *Stx2c* prophage ( $\Phi$ Stx2c). A 220 kbp inverted duplication was identified in strain Z1723<sup>Dup220</sup>. The duplicated region was flanked by repeat sequences from within prophage located at 2.2 and 2.4 Mbp relative to *OriC* and inserted into  $\Phi$ Stx2c (3.4 Mbp) bisecting the *Stx2c* prophage (Fig. 4b, c). A second 1.2 Mbp inversion was identified in strain Z1767<sup>inv1.2</sup> that also involved recombination between repeat sequences within the same prophage located at 2.2 Mbp and  $\Phi$ Stx2c (Fig. 4b, c). PFGE analysis confirmed the presence of the LCRs in Z1723<sup>Dup220</sup> and Z1767<sup>inv1.2</sup> (Fig. S6a).

## Real-time occurrence of LCRs during *in vitro* laboratory culture

We investigated if LCRs could be generated and detected in real-time following standard laboratory culture of bacteria in LB media. To increase the sensitivity of detection we applied both Oxford Nanopore Technologies (ONT) long-read sequencing and optical mapping to detect LCRs. We focused

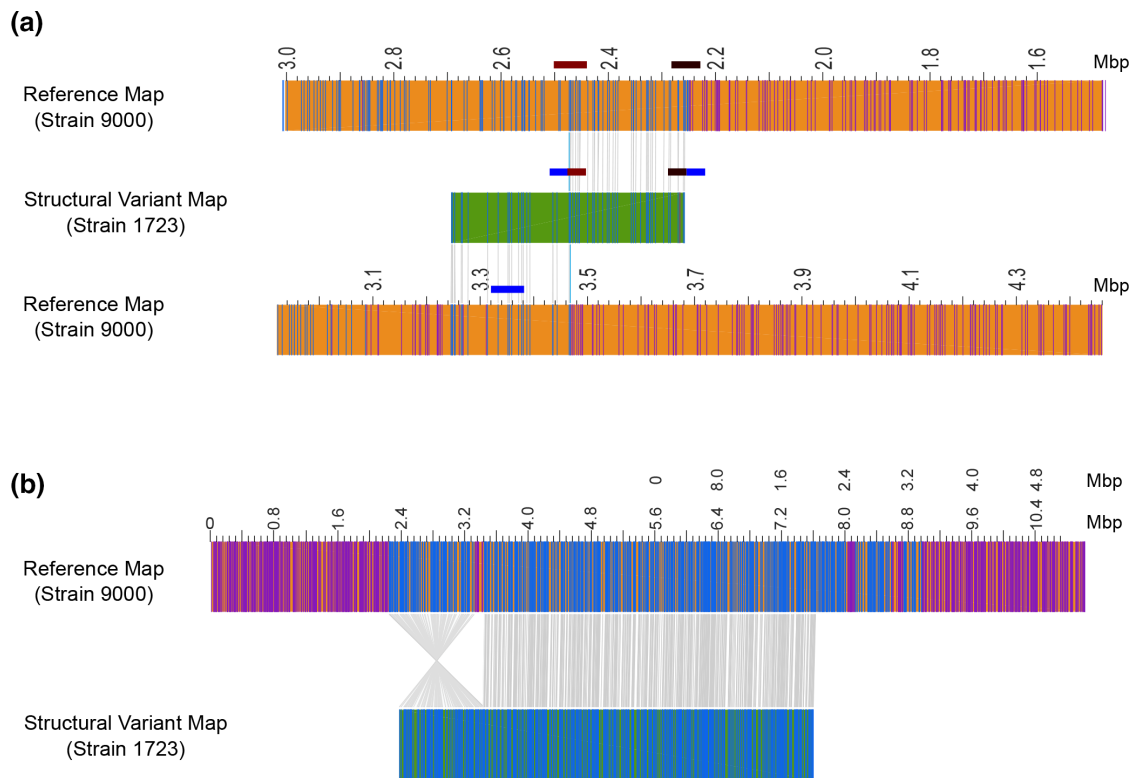


**Fig. 3.** Distinct PFGE restriction patterns of *E. coli* O157:H7 PT21/28 strains are largely accounted for by LCRs. (a) Phylogenetic distribution of lineage Ic PT21/28 strains. The strain source attribution, human (red) or bovine (blue), for each strain is indicated at each branch tip. Where available, the distinct PFGE profile types of PT21/28 strains are indicated at branch tips (coloured blocks) and those with profile C (A\_11b) are highlighted (black blocks). (b) Pairwise whole-genome comparison of ten PT21/28 strains with different PFGE profiles. Whole genomes (black lines) are centred by the replication terminus (Ter) and loci of prophage (yellow boxes), Stx prophage ( $\Phi$ Stx2c; blue and  $\Phi$ Stx2a; red) and AvrII sites (blue triangles) are shown. Direct (purple) and inverted (orange) homology at a BLAST cut-off of 10000 bp between strains are plotted. (c) PFGE profile of the ten selected PT21/28 strains following AvrII digestion. (d) *In silico* generated AvrII digestion pattern of the PacBio-generated sequences for each strain (red arrows) and bands of difference between *in silico* and AvrII digestion of strain Z892 (pink arrows) are highlighted.



**Fig. 4.** Detection of LCRs in *E. coli* O157:H7 PT21/28 strain 9000 variants analysed from cattle colonization studies. Pairwise whole-genome comparisons of strains from trial 1 (a) and trial 2 (b) are shown with direct (purple) and inverted (orange) homology at a blast cut-off of 10000 bp between strains. Whole genomes (black lines) are centred by the replication terminus (Ter) and the loci of prophage (yellow boxes) and Stx prophage ( $\Phi$ Stx2c; blue and  $\Phi$ Stx2a; red) in each strain are mapped. (c) Circos plots showing the identified 220 kbp duplication in Z1723 (left) and 1.2 Mbp inversion in Z1767 (right) relative to progenitor strain 9000. Outer ring: strain 9000 (grey) and LCR derivatives Z1723 and Z1767 (black); middle ring: loci of prophage (black) and prophage a LCR boundaries (red); inner ring: GC content.





**Fig. 5.** Optical mapping of *E. coli* O157:H7 PT21/28 strain 9000 variant Z1723. Structural variant (SV) analysis identified a 220 Kb duplication (a) and 1.2 Mb inversion (b) in the population of Z1723 relative to the reference strain 9000. The genome map (orange) of reference strain 9000 and each Z1723 structural variant (green) are shown. All restriction sites used to generate maps are shown. Paired restriction sites (blue lines) are aligned between the reference and variant maps (grey lines). Unpaired restriction sites (purple lines) outside aligned regions are also shown. The SV map containing the 220 Kb duplication has been aligned to two reference strain 9000 genome maps to demonstrate the hybrid composition of the map containing  $\Phi$ Stx2c at 3.4 Mb and an inverted 220 Kb duplicated region originating from between 2.2 and 2.4 Mb. The position of phage regions flanking the 220 kbp duplication (maroon and black blocks) and  $\Phi$ Stx2c (blue block) into which the duplication has inserted are indicated.

on analysis of strain 9000 and strain derivatives isolated from the animal colonization studies.

The wild-type parental strain 9000 was first sequenced using ONT and analysed for reads that aligned to the LCRs identified in variant strains Z1615<sup>inv1.4</sup>, Z1767<sup>inv1.2</sup> or Z1723<sup>Dup220</sup>. Aligning strain 9000 reads to the Z1615<sup>inv1.4</sup> genome, a total of five reads were found that matched the identified 1.4 Mb inversion boundary at 1.95 Mb relative to OriC and a single read that matched the inversion boundary at 3.35 Mb. These reads were approximately 2 and 0.33%, respectively, of the total reads across the same region that mapped directly to strain 9000. Similarly aligning 9000 reads to the Z1767<sup>inv1.2</sup> genome, a single read (0.4% abundance) was found that matched the 1.2 Mb inversion boundary at 2.25 Mb relative to OriC and three reads (1.2% abundance) that matched the inversion boundary within  $\Phi$ Stx2c at 3.45 Mb. No reads were found that mapped to the 220 kb duplication in Z1723<sup>Dup220</sup>.

Next we analysed strains Z1723<sup>Dup220</sup> and Z1767<sup>inv1.2</sup> using Bionano Irys optical mapping (Figs 5 and S6) to identify additional LCRs that occur during growth in LB medium. Cultures of each strain were started from single colonies and

chromosomes were extracted during late exponential phase cultures ( $OD_{600}=0.7$ ). Structural variant (SV) analysis was performed to detect all novel genome restriction maps within the cultured populations of Z1723<sup>Dup220</sup> and Z1767<sup>inv1.2</sup> that did not map directly to an *in silico* generated map of the parental strain 9000 reference genome (Fig. 5).

Optical mapping showed that both strains had mixed population structures when cultured *in vitro*. SV analysis confirmed the same 220 kb inverted duplication was present in the Z1723<sup>Dup220</sup> population that was identified by PacBio sequencing and PFGE (Fig. 5a). This hybrid structural variant mapped 5' – 3' between 2.24–2.46 Mb and 3' – 5' between 3.26–3.46 Mb to strain 9000 further confirming the presence of the inverted duplication within  $\Phi$ Stx2c at 3.4 Mb. A 1.2 Mb inversion relative to strain 9000 (Fig. 5b) was also identified in Z1723<sup>Dup220</sup>. This inversion matched the 1.2 Mb inversion seen in strain Z1767<sup>inv1.2</sup> (Fig. 4b) with boundaries in prophage located at 2.2 Mbp and 3.4 Mbp ( $\Phi$ Stx2c). PFGE analysis of two separate Z1723<sup>Dup220</sup> freezer stocks (Fig. S6a) showed that the 220 kbp inverted duplication was the dominant genome conformation present with no evidence

of secondary bands indicative of the Z1767<sup>inv1.2</sup> inversion. We therefore assume that the 1.2 Mbp inversion detected in the Z1723<sup>Dup220</sup> population by optical mapping is a minority population below the limit of detection by PFGE.

SV analysis of Z1767<sup>inv1.2</sup> identified the expected 1.2 Mbp inversion relative to strain 9000 (Fig. S7a) as determined from Pac-Bio sequencing and identified a novel 140.5 kbp inverted duplication within the cultured population (Fig. S7b). The duplicated region spanned 2.1–2.24 Mbp relative to OriC and was flanked by prophage sequence (2.2 Mbp) and an IS66 sequence located within the O-Island 48 [61]. This duplicated region also inserted in an inverted orientation within the Stx2c prophage further highlighting  $\Phi$ Stx2c as a hotspot for recombinational events leading to LCRs.

### Changes in bacterial gene expression and phenotypes associated with LCRs

Using the structural variants of strain 9000 (Z1723<sup>Dup220</sup>, Z1767<sup>inv1.2</sup>, Z1615<sup>inv1.4</sup>) generated during *in vivo* colonization we examined if the identified LCRs impacted strain gene expression and phenotypes. The global transcriptomes of strain 9000 and each structural variant strain (Z1723<sup>Dup220</sup>, Z1767<sup>inv1.2</sup>, Z1615<sup>inv1.4</sup>) were first compared by RNAseq for two growth conditions: nutrient-rich LB medium and M9 minimal medium. M9 minimal medium was chosen as the transcriptional response of *E. coli* O157:H7 in M9 was previously shown to be similar to that observed in cattle faeces compared with LB [62]. PCA analysis showed there was little discernible difference between the transcriptomes of each strain when cultured in LB (Fig. S8a) however the transcriptome of strain Z1723<sup>Dup220</sup>, containing a 220kbp inverted duplication, was distinct from strain 9000 and the other variants in M9 (Fig. S8b). Differential changes in gene expression were modest (Table S2) although a gene dosage effect was apparent across the region of duplication with an increase in expression observed for 66 of the duplicated genes when mapped to the genome of WT strain 9000 (Fig. 6a). There was also a marked effect on the expression of genes within the Stx2a prophage ( $\Phi$ Stx2a) rather than the Stx2c prophage ( $\Phi$ Stx2c) into which the 220 kbp duplication had inserted (Fig. 6a).

As Stx2 toxin is the primary virulence factor of *E. coli* O157:H7 strains leading to HUS, we tested if the observed differential transcription within  $\Phi$ Stx2a in Z1723<sup>Dup220</sup> affected Stx2a expression, production and activity compared with other structural variants. For each phenotype Z1723<sup>Dup220</sup> was compared with trial 2 variants Z1766 and Z1767<sup>inv1.2</sup>. Strains 9000 and Z1615<sup>inv1.4</sup> were excluded due to the previously documented [52] inactivation of the *stx2a* gene by an IS element, IS629. Expression of *stx2a* was increased in Z1723<sup>Dup220</sup> compared to both Z1766 and Z1767<sup>inv1.2</sup> (Fig. 6b) and this manifested as a significant increase in total Stx2 toxin (Fig. 6c) and cytotoxic killing of Stx2 susceptible Vero cells (Fig. 6d).

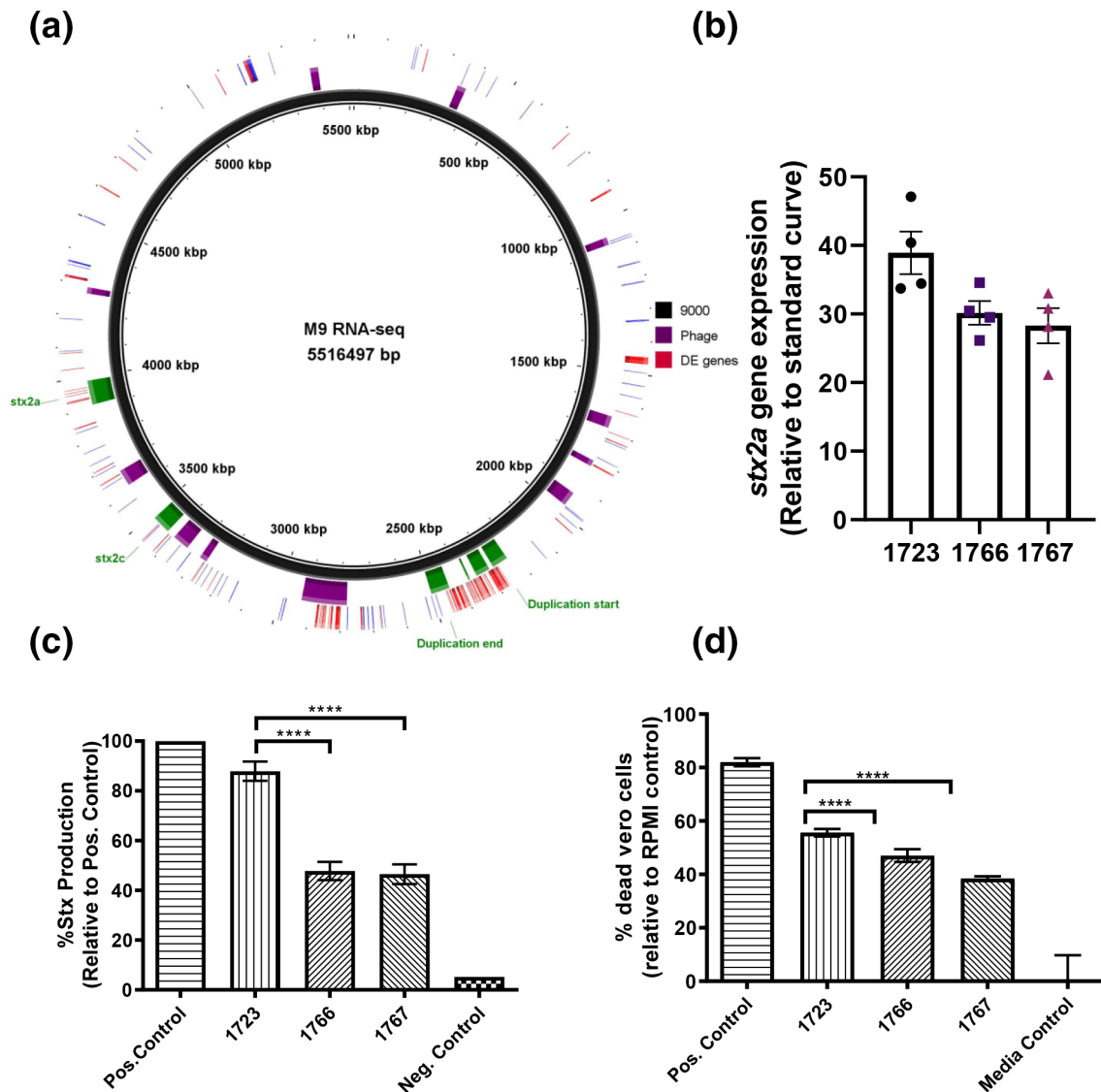
We have previously shown that lysogeny with Stx2 prophages negatively regulates the LEE type III secretion system (T3SS)

[56] and demonstrated that a large duplication may have influenced the fitness of two closely related outbreak strains [63]. As the 220 kbp duplication in Z1723<sup>Dup220</sup> interrupted  $\Phi$ Stx2c and increased expression of  $\Phi$ Stx2a genes we examined T3S and assessed the competitive fitness for strains Z1723<sup>Dup220</sup>, Z1767<sup>inv1.2</sup> and Z1766 (Fig. 7). Transcriptional *gfp* fusions to the LEE master regulator, *ler*, and LEE4 encoded *sepL* were introduced into each strain and expression was monitored in MEM-HEPES medium (OD<sub>600</sub>=0.8). Expression of both *ler* and *sepL* was decreased in Z1723<sup>Dup220</sup> compared to Z1766 and Z1767<sup>inv1.2</sup> (Fig. 7a). There was also a marked difference in the levels of the T3S secreted protein, EspD, which could not be detected in the culture supernatant of Z1723<sup>Dup220</sup> (Fig. 7b).

The competitive fitness of strains from trial 1 (9000 and Z1615<sup>inv1.4</sup>) and trial 2 (Z1723<sup>Dup220</sup>, Z1767<sup>inv1.2</sup> and Z1766) was assessed by paired co-culturing in M9 media. In M9 media Z1723<sup>Dup220</sup> significantly outcompeted the structural variants Z1766 and Z1767<sup>inv1.2</sup> as mean fitness indices (f.i.) of 0.89 and 0.93 were recorded, respectively (Fig. 7c) compared with control, f.i.=1. No significant difference in fitness was observed between trial 1 strains in M9 (Fig. 7c). Finally, we measured the motility of strains 9000 and each structural variant on tryptone swarm plates (Fig. 7d). The *flg* and *fli* motility gene operons are naturally located on opposing sides of the Z1615<sup>inv1.4</sup> and Z1767<sup>inv1.2</sup> inversions and in different chromosomal domains. Each inversion brings these operons proximal to each other and we determined if this positional shift impacted motility. For strains isolated from calf trial 2 no difference in motility between Z1723<sup>Dup220</sup> and Z1767<sup>inv1.2</sup> was observed however Z1766 was significantly more motile than both variants. Z1615<sup>inv1.4</sup> from calf trial 1 was also significantly more motile than WT strain 9000. These data provide evidence that LCRs can impact important *E. coli* O157:H7 phenotypes involved in host colonization and disease.

## DISCUSSION

This study describes a systematic comparison of genome structure across the main lineages of the zoonotic pathogen, *E. coli* O157:H7. As originally demonstrated for a single strain, EDL933 [19], LCRs occur across the main lineages of this serogroup and are bounded by specific prophages clustered towards the terminus of the genome. Large prophage homologous repeats (>5000 bp) were identified at the boundaries of LCRs, which provide ample sequence substrate for recombination leading to duplications and inversions. While it is likely that LCRs are mediated by RecABCD recombination, *E. coli* O157:H7 strains also carry multiple  $\lambda$ -like phage, including Stx phage, many of which encode their own Rad52-like recombinase enzymes such as Red $\beta$  [64–66]. Whether the formation of LCRs in *E. coli* O157:H7 strains is host or phage mediated (or both) still needs to be determined. Irrespective of the recombination system involved, the generation of LCRs would require a double-strand break (DSB) in one or more of the phage at their boundaries. It is interesting to speculate that double-strand breaks (DSBs) within phage are a primary driver of recombinational repair in bacteria via the SOS



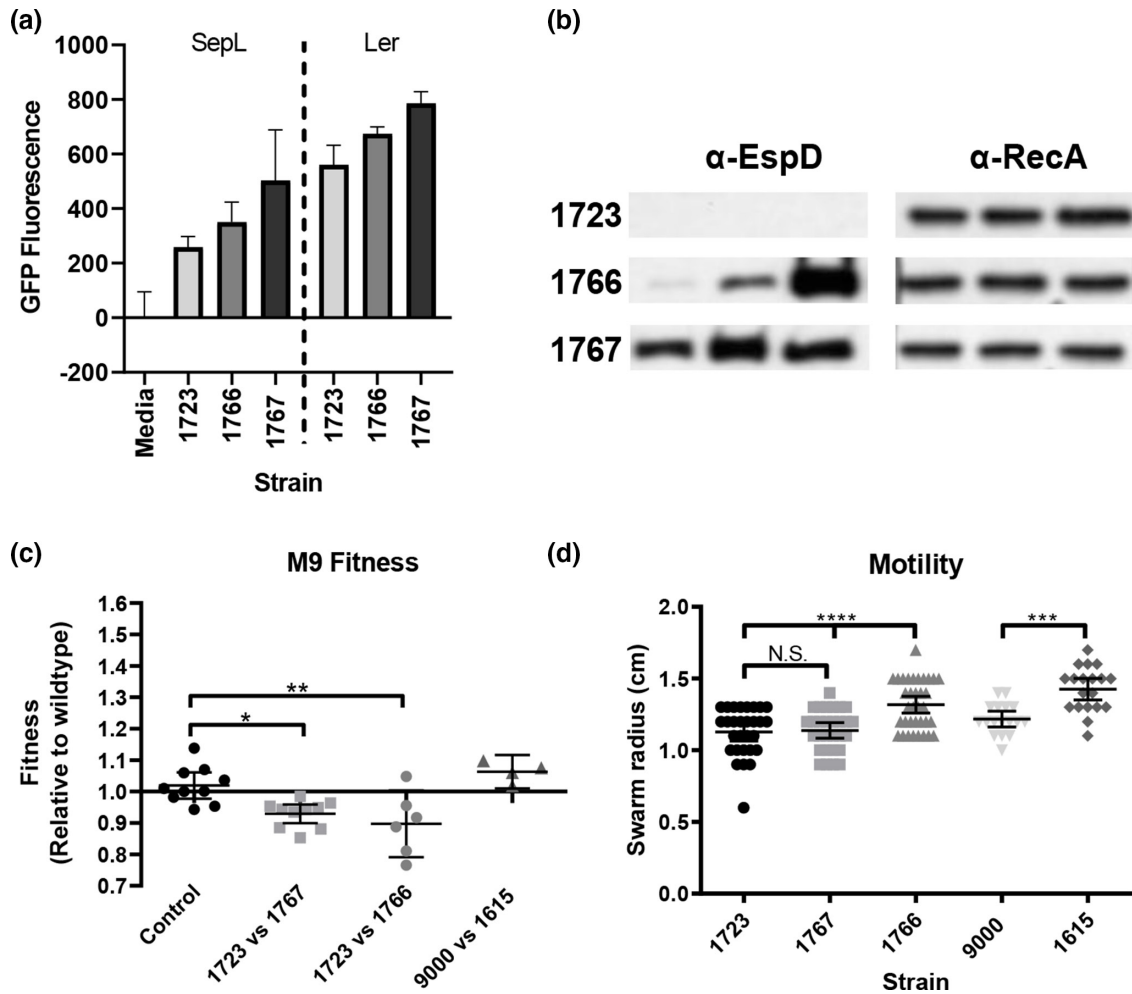
**Fig. 6.** Shiga-toxin expression, production and toxicity of strain 9000 structural variants. The chromosomal location of all differentially expressed genes in Z1723 (upregulated, red; downregulated, blue) are mapped to the reference strain 9000 genome (a). Prophage (purple blocks), the 220 kb region of duplication (green),  $\Phi$ Stx2a and  $\Phi$ Stx2c regions are highlighted. Expression of *stx2a* (b) total Stx toxin production (c) and Vero cell toxicity Stx (d) was measured for trial 2 strains Z1723, Z1766 and Z1767 in M9 media. Mean values  $\pm$  SEM of four biological replicates ( $n=4$ ) are shown for each assay. \* $P \leq 0.05$ ; \*\* $P \leq 0.01$ ; \*\*\* $P \leq 0.001$ ; \*\*\*\* $P \leq 0.0001$ .

response that is also required for prophage-based expression of Stx [67]. Strains of *E. coli* O157:H7 PT21/28 constitutively express Stx2 [52] and therefore the rate of occurrence of DSBs, RecA-mediated Stx expression and LCR formation may be interconnected.

The Stx2c-encoding prophage was shown to be present across the different *E. coli* O157:H7 lineages without much variation compared to Stx2a-encoding prophages [9, 11]. In the present study the Stx2c prophage is a primary site of many of the LCRs. One structural variant of PT21/28 strain 9000 was a duplication from one side of the genome that inserted into the Stx2c terminase gene region on the other side of the genome. This was of particular interest as this large region

of duplicated homology could stimulate inversions and also recombination resolving back to the original confirmation. A similar duplication has been sequenced in two closely related strains associated with sequential *E. coli* O157:H7 outbreaks at a single restaurant [19].

The ONT long-read sequencing and optical mapping results provide evidence that LCRs are continuously generated. The estimation from the ONT long-read sequencing of strain 9000 was between 1–2% of the population when cultured in LB. Specific LCRs were also detected during animal colonization but it remains to be demonstrated if specific variants are being selected under *in vivo* conditions, i.e. that such variants offer any fitness advantage in real environments.



**Fig. 7.** Type III secretion, competitive fitness and motility phenotypes of strain 9000 structural variants. (a) Expression of the LEE master regulator *ler* and LEE4 chaperone *sepL* was measured by Gfp reporter fusions ( $n=3$ ). (b) Detection of the LEE effector EspD in the culture supernatants of each strain by Western blot ( $n=3$ ). Corresponding cellular RecA levels were used as a control. (c) Competitive fitness of strains after 24 h co-culturing in M9 media ( $n=6$ ). (d) Motility of strains after 6 h on Tryptone swarm plates ( $n=20$ ). Mean values  $\pm$  SEM for  $n$  biological replicates are shown for each assay. \* $P \leq 0.05$ ; \*\* $P \leq 0.01$ ; \*\*\* $P \leq 0.001$ ; \*\*\*\* $P \leq 0.0001$ .

Currently, there is no cost-effective way to quantify the proportions of the conformational variants from complex animal samples as outgrowth of the populations *in vitro* is required to work with the long-read sequencing analyses or optical mapping. We show that LCRs likely account for the majority of observed PFGE variation among PT21/28 bovine isolates in the UK. One exception was strain Z892 in which the observed PFGE profile differed from the *in silico* digestion profile. As LCRs occur during routine culture it is possible that the observed profile change resulted from an LCR variant coming to dominance within the population. Phage loss or movement could also drive such changes in PFGE profiles as shown for *E. coli* O157:H7 strain Sakai [68] and IS element movements and SNPs could also be involved. Re-sequencing strain Z892 would be required to determine the exact nature of this difference. Extensive surveys of *E. coli* O157:H7 in cattle herds [69, 70] has determined that while the majority of isolates in any specific herd exhibit

the same PFGE pattern, there are isolates with different profiles yet the same phage type (PT) [71, 72], indicative that LCRs may also be occurring in strains present within a herd. A recent study of persistent lineage I strains isolated on a single farm also demonstrated that a 47.7 kbp deletion was a significant genomic difference between two of the strains [73]. There has been one previous report of multiple deletions occurring during *E. coli* O157:H7 colonization of cattle, generating multiple PFGE types [74].

LCRs have been observed in a number of bacterial genera, including *Campylobacter*, *Yersinia*, *Staphylococcus* and *Salmonella* [22, 75–78] and duplications leading to rearrangements in gene order have recently been proposed to play a key role in niche adaptation [79]. For inversions the gene content and copy number is maintained but the prophage boundaries do change in composition and this could have an impact on prophage gene expression or the regulatory networks



that they are part of [1]. A clear example of this was shown for *Campylobacter* where in one orientation the inversion completes an active prophage and in turn that provides resistance to certain infecting phages [75]. For *E. coli* O157:H7 strain 9000 structural variants we measured a number of expression and phenotype changes, including motility and growth rate for variants with inversions. The most obvious differences were present in the variant with a 220 kbp duplication. This included an increase in Stx expression, production and toxicity and a reduction in type 3 secretion. Our previous research has shown that Stx2a prophage integration into different *E. coli* backgrounds led to a repression of T3S, potentially via the CII protein [56]. Such cross-regulation would offer one pathway resulting in the concomitant reduction in T3S in the strain with the duplication. Chromosomal inversions involving the Ter region that lead to replicore imbalance can stall or stop replication forks and induce SOS [80]. Due to the spatial distribution of prophages involved in LCRs, the main large inversions we have identified generally do not generate major changes in replicore size. However even minor changes could impact growth rate and phenotypes as seen with the LCR specific phenotypes identified in this study affecting virulence gene expression, fitness and motility.

By definition, phenotypic heterogeneity is the occurrence of individuals within a genetically identical population that stochastically develop phenotypes of varying fitness within a homogenous environment [81–84]. With the work presented here and that in other genera, it is evident that genome structural variants can be a way to generate phenotypic heterogeneity in a clonal bacterial population. More studies are required to determine if LCR sub-populations provide fitness advantages under different growth conditions, which would mean that such recombinational capacity would represent an important survival strategy.

#### Funding information

We acknowledge funding from a Food Standards Scotland/Food Standards Agency research grant (FSS101055) along with a BBSRC research grant BB/P02095X/1. We also acknowledge the support from an institute core strategic grant to the Roslin Institute: BBS/E/D/20002173.

#### Acknowledgements

The authors would like to thank Sandy Fryda-Bradley and the USMARC core sequencing facility for excellent technical assistance. The mention of a trade name, proprietary product, or specific equipment does not constitute a guarantee or warranty by the USDA and does not imply approval to the exclusion of other products that might be suitable. The USDA is an equal opportunity employer and provider.

#### Conflicts of interest

The authors declare that there are no conflicts of interest

#### References

- Brussow H, Canchaya C, Hardt WD. Phages and the evolution of bacterial pathogens: from genomic rearrangements to lysogenic conversion. *Microbiol Mol Biol Rev* 2004;68:560–602.
- Taylor VL, Fitzpatrick AD, Islam Z, Maxwell KL. The diverse impacts of phage morons on bacterial fitness and virulence. *Adv Virus Res* 2019;103:S0065-3527(18)30043-5:1–31..
- Argov T, Azulay G, Pasechnek A, Stadnyuk O, Ran-Sapir S. Temperate bacteriophages as regulators of host behavior. *Curr Opin Microbiol* 2017;38:81–87.
- Bondy-Denomy J, Qian J, Westra ER, Buckling A, Guttman DS. Prophages mediate defense against phage infection through diverse mechanisms. *ISME J* 2016;10:2854–2866.
- Davies EV, Winstanley C, Fothergill JL, James CE. The role of temperate bacteriophages in bacterial infection. *FEMS Microbiol Lett* 2016;363:fnw015.
- Wang X, Kim Y, Ma Q, Hong SH, Pokusaeva K. Cryptic prophages help bacteria cope with adverse environments. *Nat Commun* 2010;1:147.
- Tree JJ, Granneman S, McAteer SP, Tollervey D, Gally DL. Identification of bacteriophage-encoded anti-sRNAs in pathogenic *Escherichia coli*. *Mol Cell* 2014;55:199–213.
- Ferens WA, Hovde CJ. *Escherichia coli* O157:H7: animal reservoir and sources of human infection. *Foodborne Pathog Dis* 2011;8:465–487.
- Shaaban SC, McAteer LA, Jenkins SP, Dallman C, Bono JL, et al. Evolution of a zoonotic pathogen: investigating prophage diversity in enterohaemorrhagic *Escherichia coli* O157 by long-read sequencing. *Microb Genom* 2016.
- Perna NT, Plunkett G 3rd, Burland V, Mau B, Glasner JD, et al. Genome sequence of enterohaemorrhagic *Escherichia coli* O157:H7. *Nature* 2001;409:529–533.
- Dallman TJ, Ashton PM, Byrne L, Perry NT, Petrovska L, et al. Applying phylogenomics to understand the emergence of Shiga-toxin-producing *Escherichia coli* O157:H7 strains causing severe human disease in the UK. *Microb Genom* 2015;1:e000029.
- Heiman KE, Mody RK, Johnson SD, Griffin PM, Gould LH. *Escherichia coli* O157 Outbreaks in the United States, 2003–2012. *Emerging Infect Dis* 2015;21:1293–1301.
- Karmali MA. Host and pathogen determinants of verocytotoxin-producing *Escherichia coli*-associated hemolytic uremic syndrome. *Kidney Int Suppl* 2009;S4–7.
- Obrig TG, Karpman D. Shiga toxin pathogenesis: kidney complications and renal failure. *Curr Top Microbiol Immunol* 2012;357:105–136.
- Brandal LT, Wester AL, Lange H, Lobersli I, Lindstedt BA. Shiga toxin-producing *Escherichia coli* infections in Norway, 1992–2012: characterization of isolates and identification of risk factors for haemolytic uremic syndrome. *BMC Infect Dis* 2015;15:324.
- Iyoda S, Manning SD, Seto K, Kimata K, Isobe J. Phylogenetic Clades 6 and 8 of Enterohemorrhagic *Escherichia coli* O157:H7 with particular stx subtypes are more frequently found in isolates from hemolytic uremic syndrome patients than from asymptomatic carriers. *Open Forum Infect Dis* 2014;1:ofu061.
- Buvens G, De Gheldre Y, Dediste A, de Moreau AI, Mascart G. Incidence and virulence determinants of verocytotoxin-producing *Escherichia coli* infections in the Brussels-Capital Region, Belgium, in 2008–2010. *J Clin Microbiol* 2012;50:1336–1345.
- Fuller CA, Pellino CA, Flagler MJ, Strasser JE, Weiss AA. Shiga toxin subtypes display dramatic differences in potency. *Infect Immun* 2011;79:1329–1337.
- Iguchi A, Iyoda S, Terajima J, Watanabe H, Osawa R. Spontaneous recombination between homologous prophage regions causes large-scale inversions within the *Escherichia coli* O157:H7 chromosome. *Gene* 2006;372:199–207.
- Kotewicz ML, Jackson SA, LeClerc JE, Cebula TA. Optical maps distinguish individual strains of *Escherichia coli* O157: H7. *Microbiology (Reading)* 2007;153:1720–1733.
- Periwal V, Scaria V. Insights into structural variations and genome rearrangements in prokaryotic genomes. *Bioinformatics* 2015;31:1–9.
- Hughes D. Evaluating genome dynamics: the constraints on rearrangements within bacterial genomes. *Genome Biol* 2000;1:REVIEWS0006.

23. Clawson ML, Keen JE, Smith TP, Durso LM, McDaneld TG. Phylogenetic classification of *Escherichia coli* O157:H7 strains of human and bovine origin using a novel set of nucleotide polymorphisms. *Genome Biol* 2009;10:R56.
24. Koren S, Harhay GP, Smith TP, Bono JL, Harhay DM. Reducing assembly complexity of microbial genomes with single-molecule sequencing. *Genome Biol* 2013;14:R101.
25. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* 2013;10:563–569.
26. Stewart AC, Osborne B, Read TD. DIYA: a bacterial annotation pipeline for any genomics lab. *Bioinformatics* 2009;25:962–963.
27. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, et al. Pilon: An integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS one* 2014;9:e112963.
28. Luo H, Zhang CT, Gao F. Ori-Finder 2, an integrated tool to predict replication origins in the archaeal genomes. *Front Microbiol* 2014;5:482.
29. Arndt D, Grant JR, Marcu A, Sajed T, Pon A. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res* 2016;44:W16–21.
30. Green MR, Sambrook J. Isolation and quantification of DNA. *Cold Spring Harb Protoc* 2018;2018:pdb.top093336.
31. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;34:3094–3100.
32. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 2009;25:2078–2079.
33. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 2010;26:841–842.
34. Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer. *Bioinformatics* 2011;27:1009–1010.
35. Okonechnikov K, Golosova O, Fursov M, team U. Unipro UGENE: a unified bioinformatics toolkit. *Bioinformatics* 2012;28:1166–1167.
36. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol* 1990;215:403–410.
37. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R. Circos: an information aesthetic for comparative genomics. *Genome Res* 2009;19:1639–1645.
38. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* 2015;31:3691–3693.
39. Price MN, Dehal PS, Arkin AP. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS one* 2010;5:e9490.
40. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments. *Nucleic Acids Res* 2019;47:W256–W259.
41. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinformatics* 2010;26:589–595.
42. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;20:1297–1303.
43. Dallman T, Ashton P, Schafer U, Jironkin A, Painset A. SnapperDB: a database solution for routine sequencing analysis of bacterial isolates. *Bioinformatics* 2018;34:3028–3029.
44. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res* 2015;43:e15.
45. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol* 2020;37:1530–1534.
46. Ribot EM, Fair MA, Gautom R, Cameron DN, Hunter SB. Standardization of pulsed-field gel electrophoresis protocols for the subtyping of *Escherichia coli* O157:H7, *Salmonella*, and *Shigella* for PulseNet. *Foodborne Pathog Dis* 2006;3:59–67.
47. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29:15–21.
48. Hyatt D, Chen GL, Locascio PF, Land ML, Larimer FW. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 2010;11:119.
49. Nikolayeva O, Robinson MD. edgeR for differential RNA-seq and ChIP-seq analysis: an application to stem cell biology. *Methods Mol Biol* 2014;1150:45–79.
50. Dykhuizen DE. Experimental studies of natural selection in bacteria. *Annu Rev Ecol Syst* 1990;21:373–398.
51. Lenski RE. Quantifying fitness and gene stability in microorganisms. *Biotechnology* 1991;15:173–192.
52. Fitzgerald SF, Beckett AE, Palarea-Albaladejo J, McAteer S, Shaaban S, et al. Shiga toxin sub-type 2a increases the efficiency of *Escherichia coli* O157 transmission between animals and restricts epithelial regeneration in bovine enteroids PLoS pathogens. *Research* 2019.
53. Fernandez-Brando RJ, Yamaguchi N, Tahoun A, McAteer SP, Gillespie T. Type III secretion-dependent sensitivity of *Escherichia coli* O157 to specific ketolides. *Antimicrob Agents Chemother* 2016;60:459–470.
54. Wang D, Roe AJ, McAteer S, Shipston MJ, Gally DL. Hierarchical type III secretion of translocators and effectors from *Escherichia coli* O157:H7 requires the carboxy terminus of SepL that binds to Tir. *Mol Microbiol* 2008;69:1499–1512.
55. Holden N, Totsika M, Dixon L, Catherwood K, Gally DL. Regulation of P-fimbrial phase variation frequencies in *Escherichia coli* CFT073. *Infect Immun* 2007;75:3325–3334.
56. Xu X, McAteer SP, Tree JJ, Shaw DJ, Wolfson EB. Lysogeny with Shiga toxin 2-encoding bacteriophages represses type III secretion in enterohemorrhagic *Escherichia coli*. *PLoS Pathog* 2012;8:e1002672.
57. Doyle M, Fookes M, Ivens A, Mangan MW, Wain J. An H-NS-like stealth protein aids horizontal DNA transmission in bacteria. *Science* 2007;315:251–252.
58. Darmon E, Leach DR. Bacterial genome instability. *Microbiol Mol Biol Rev* 2014;78:1–39.
59. Ooka T, Ogura Y, Asadulghani M, Ohnishi M, Nakayama K. Inference of the impact of insertion sequence (IS) elements on bacterial genome diversification through analysis of small-size structural polymorphisms in *Escherichia coli* O157 genomes. *Genome Res* 2009;19:1809–1816.
60. Corbishley A, Ahmad NI, Hughes K, Hutchings MR, McAteer SP. Strain-dependent cellular immune responses in cattle following *Escherichia coli* O157:H7 colonization. *Infect Immun* 2014;82:5117–5131.
61. Taylor DE, Rooker M, Keelan M, Ng LK, Martin I. Genomic variability of O islands encoding tellurite resistance in enterohemorrhagic *Escherichia coli* O157:H7 isolates. *J Bacteriol* 2002;184:4690–4698.
62. Landstorfer R, Simon S, Schober S, Keim D, Scherer S. Comparison of strand-specific transcriptomes of enterohemorrhagic *Escherichia coli* O157:H7 EDL933 (EHEC) under eleven different environmental conditions including radish sprouts and cattle feces. *BMC genomics* 2014;15:353.
63. Cowley LA, Dallman TJ, Fitzgerald S, Irvine N, Rooney PJ, et al. Short-term evolution of Shiga toxin-producing *Escherichia coli* O157:H7 between two food-borne outbreaks. *Microb Genom* 2016;2:e000084.
64. Bobay LM, Touchon M, Rocha EP. Manipulating or superseding host recombination functions: a dilemma that shapes phage evolvability. *PLoS Genet* 2013;9:e1003825.
65. De Paepe M, Hutinet G, Son O, Amarir-Bouhram J, Schbath S. Temperate phages acquire DNA from defective prophages by relaxed homologous recombination: the role of Rad52-like recombinases. *PLoS Genet* 2014;10:e1004181.
66. Asadulghani M, Ogura Y, Ooka T, Itoh T, Sawaguchi A. The defective prophage pool of *Escherichia coli* O157: prophage-prophage interactions potentiate horizontal transfer of virulence determinants. *PLoS Pathog* 2009;5:e1000408.

67. Schmidt H. Shiga-toxin-converting bacteriophages. *Res Microbiol* 2001;152:687–695.
68. Chen C, Lewis CR, Goswami K, Roberts EL, DebRoy C. Identification and characterization of spontaneous deletions within the Sp11-Sp12 prophage region of *Escherichia coli* O157:H7 Sakai. *Appl Environ Microbiol* 2013;79:1934–1941.
69. Henry MK, Tongue SC, Evans J, Webster C, MC KI. British *Escherichia coli* O157 in Cattle Study (BECS): to determine the prevalence of *E. coli* O157 in herds with cattle destined for the food chain. *Epidemiol Infect* 2017;145:3168–3179.
70. Matthews L, McKendrick IJ, Ternent H, Gunn GJ, Synge B. Super-shedding cattle and the transmission dynamics of *Escherichia coli* O157. *Epidemiol Infect* 2006;134:131–142.
71. Pearce MC, Jenkins C, Vali L, Smith AW, Knight HI. Temporal shedding patterns and virulence factors of *Escherichia coli* serogroups O26, O103, O111, O145, and O157 in a cohort of beef calves and their dams. *Appl Environ Microbiol* 2004;70:1708–1716.
72. Vali L, Wisely KA, Pearce MC, Turner EJ, Knight HI. High-level genotypic variation and antibiotic sensitivity among *Escherichia coli* O157 strains isolated from two Scottish beef cattle farms. *Appl Environ Microbiol* 2004;70:5947–5954.
73. Stanton E, Wahlig TA, Park D, Kaspar CW. Chronological set of *E. coli* O157:H7 bovine strains establishes a role for repeat sequences and mobile genetic elements in genome diversification. *BMC genomics* 2020;21:562.
74. Yoshii N, Ogura Y, Hayashi T, Ajiro T, Sameshima T. Pulsed-field gel electrophoresis profile changes resulting from spontaneous chromosomal deletions in enterohemorrhagic *Escherichia coli* O157:H7 during passage in cattle. *Appl Environ Microbiol* 2009;75:5719–5726.
75. Scott AE, Timms AR, Connerton PL, Loc Carrillo C, Adzfa Radzum K. Genome dynamics of *Campylobacter jejuni* in response to bacteriophage predation. *PLoS Pathog* 2007;3:e119.
76. Darling AE, Miklos I, Ragan MA. Dynamics of genome rearrangement in bacterial populations. *PLoS Genet* 2008;4:e1000128.
77. Guerillot R, Kostoulas X, Donovan L, Li L, Carter GP. Unstable chromosome rearrangements in *Staphylococcus aureus* cause phenotype switching associated with persistent infections. *Proc Natl Acad Sci U S A* 2019;116:20135–20140.
78. Sun S, Ke R, Hughes D, Nilsson M, Andersson DI. Genome-wide detection of spontaneous chromosomal rearrangements in bacteria. *PLoS one* 2012;7:e42639.
79. Brandis G, Hughes D. The SNAP hypothesis: Chromosomal rearrangements could emerge from positive selection during niche adaptation. *PLoS Genet* 2020;16:e1008615.
80. Lesterlin C, Pages C, Dubarry N, Dasgupta S, Cornet F. Asymmetry of chromosome replicores renders the DNA translocase activity of FTSK essential for cell division and cell shape maintenance in *Escherichia coli*. *PLoS Genet* 2008;4:e1000288.
81. Ackermann M. A functional perspective on phenotypic heterogeneity in microorganisms. *Nat Rev Microbiol* 2015;13:497–508.
82. Martins BM, Locke JC. Microbial individuality: how single-cell heterogeneity enables population level strategies. *Curr Opin Microbiol* 2015;24:S1369–5274(15)00004–1:104–112.:
83. Grimbergen AJ, Siebring J, Solopova A, Kuipers OP. Microbial bet-hedging: the power of being different. *Curr Opin Microbiol* 2015;25:S1369–5274(15)00052–1:67–72.:
84. Zhang Z, Du C, de Barys F, Liem M, Liakopoulos A, et al. Antibiotic production in *Streptomyces* is organized by a division of labor through terminal genomic differentiation. *Sci Adv* 2020;6:eaay5781.

### Five reasons to publish your next article with a Microbiology Society journal

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

Find out more and submit your article at [microbiologyresearch.org](https://microbiologyresearch.org).