



Swansea University  
Prifysgol Abertawe



## Cronfa - Swansea University Open Access Repository

---

This is an author produced version of a paper published in:

*Molecular Ecology*

Cronfa URL for this paper:

<http://cronfa.swan.ac.uk/Record/cronfa33686>

---

### **Paper:**

Pascoe, B., Méric, G., Yahara, K., Wimalarathna, H., Murray, S., Hitchings, M., Sproston, E., Carrillo, C., Taboada, E., et. al. (2017). Local genes for local bacteria: evidence of allopatry in the genomes of transatlantic *Campylobacter* populations. *Molecular Ecology*

<http://dx.doi.org/10.1111/mec.14176>

---

This item is brought to you by Swansea University. Any person downloading material is agreeing to abide by the terms of the repository licence. Copies of full text items may be used or reproduced in any format or medium, without prior permission for personal research or study, educational or non-commercial purposes only. The copyright for any work remains with the original author unless otherwise specified. The full-text must not be sold in any format or medium without the formal permission of the copyright holder.

Permission for multiple reproductions should be obtained from the original author.

Authors are personally responsible for adhering to copyright and publisher restrictions when uploading content to the repository.

<http://www.swansea.ac.uk/iss/researchsupport/cronfa-support/>

1 **Local genes for local bacteria: evidence of allopatry in the genomes of transatlantic**  
2 ***Campylobacter* populations**

3

4 Ben Pascoe<sup>1,2</sup>, Guillaume Méric<sup>1</sup>, Koji Yahara<sup>3,4</sup>, Helen Wimalalathna<sup>5</sup>, Susan Murray<sup>4</sup>,  
5 Matthew D. Hitchings<sup>4</sup>, Emma L. Sproston<sup>6</sup>, Catherine D. Carrillo<sup>7</sup>, Eduardo N. Taboada<sup>8</sup>,  
6 Kerry K. Cooper<sup>9</sup>, Steven Huynh<sup>9</sup>, William G. Miller<sup>9</sup>, Alison J. Cody<sup>5</sup>, Keith A. Jolley<sup>5</sup>,  
7 Martin M. J. Maiden<sup>5, 10</sup>, Noel D. McCarthy<sup>5,10, 11</sup>, Xavier Didelot<sup>11</sup>, Craig T. Parker<sup>9</sup> and  
8 Samuel K. Sheppard<sup>1,2,5#</sup>

9

10 <sup>1</sup>The Milner Centre for Evolution, Department of Biology and Biochemistry, Bath University,  
11 Claverton Down, Bath, BA2 7AY, UK; <sup>2</sup>MRC CLIMB Consortium, Bath University,  
12 Claverton Down, Bath, BA2 7AY, UK; <sup>3</sup>Biostatistics Center, Kurume University, Kurume,  
13 Fukuoka 830-0011, Japan; <sup>4</sup>Swansea University Medical School, Swansea University,  
14 Singleton Park, Swansea, SA2 8PP; <sup>5</sup>Department of Zoology, University of Oxford, South  
15 Parks Road, Oxford, OX1 3PS, UK; <sup>6</sup>Bureau of Microbial Hazards, Health Canada, 251 Sir  
16 Frederick Banting Driveway, Ottawa, ON K1A 0K9, Canada; <sup>7</sup>Canadian Food Inspection  
17 Agency, 960 Carling Avenue, Ottawa, ON K1A 0Y9, Canada; <sup>8</sup>National Microbiology  
18 Laboratory at Lethbridge, Public Health Agency of Canada, PO Box 640, Township Rd. 9-1,  
19 Lethbridge, AB T1J 3Z4, Canada; <sup>9</sup>Produce Safety and Microbiology Research Unit,  
20 Agricultural Research Service, US Department of Agriculture, Albany, California,  
21 USA; <sup>10</sup>NIHR HPRU in Gastrointestinal Infections <sup>11</sup>University of Warwick, Coventry, CV4  
22 7AL, UK; <sup>11</sup>Department of Infectious Disease Epidemiology, Imperial College London,  
23 London, W2 1PG, UK

24

25 **#Address correspondence to:** Professor Samuel K. Sheppard, The Milner Centre for  
26 Evolution, Department of Biology and Biochemistry, Bath University, Claverton Down,  
27 Bath, BA2 7AY, UK. Telephone: +44(0)1225 385046; Fax: +44(0)1225 386779; Email:

28 [S.K.Sheppard@bath.ac.uk](mailto:S.K.Sheppard@bath.ac.uk)

29

30 **Running title:** *Campylobacter* biogeography

31 **Keywords:** Allopatry; *Campylobacter*; Genomics; Source attribution; Recombination;

32 Phylogeny

33

34 **Abstract**

35 The genetic structure of bacterial populations can be related to geographical locations of  
36 isolation. In some species, there is a strong correlation between geographical distances and  
37 genetic distances, which can be caused by different evolutionary mechanisms. Patterns of  
38 ancient admixture in *Helicobacter pylori* can be reconstructed in concordance with past  
39 human migration, whereas in *Mycobacterium tuberculosis* it is the lack of recombination that  
40 causes allopatric clusters. In *Campylobacter*, source attribution based on genomic data has  
41 been successful in distinguishing the infected host species, but not geographical origin. We  
42 investigate biogeographical signals in highly recombining genes to determine the extent of  
43 clustering between genomes from geographically distinct *Campylobacter* populations. Whole  
44 genome sequences from 294 *Campylobacter* isolates from North America and the UK were  
45 analysed. Isolates from within the same country shared more recently recombined DNA than  
46 isolates from different countries. Using 15 UK/American pairs of isolates that shared  
47 ancestors, we identify regions that have frequently and recently recombined to test their  
48 correlation with geographical origin. The seven genes that demonstrated the greatest  
49 clustering by geography were used in an attribution model to infer geographical origins. A  
50 further 383 UK clinical isolates were used to detect signals of foreign travel. Patient records  
51 indicated that 46 cases had travelled abroad less than two weeks prior to sampling, and 34  
52 (74%) of those *Campylobacter* genomes, were deemed to be from a non-UK origin.  
53 Detection of signals of biogeographical differences in *Campylobacter* genomes will  
54 contribute to improved source attribution of clinical *Campylobacter* infection and inform  
55 intervention strategies to reduce campylobacteriosis.

56

57 **Introduction**

58

59 Geographical structuring is well documented in bacteria and analysing genetic variation  
60 among isolates can provide information about the global spread of important pathogens. For  
61 example, after spreading with Neolithic human hosts (Comas et al., 2013), lineages of  
62 *Mycobacterium tuberculosis* populations can be classified into geographical groups based  
63 upon local genetic diversification of DNA sequences (Achtman, 2008, Gagneux and Small,  
64 2007). Phylogeographic structuring has also been observed in the human stomach bacterium  
65 *Helicobacter pylori*, where a rapidly evolving genome, with high levels of horizontal gene  
66 transfer (HGT), allows the reconstruction of recent human migrations to the extent that  
67 genetic admixture among the bacteria reflects interactions among human populations (Falush  
68 et al., 2003, Moodley et al., 2009).

69

70 Tuberculosis and *H. pylori* are primarily human pathogens, but for *Campylobacter*, animals  
71 are the principal reservoir for human infection. International trade, particularly in agricultural  
72 animals such as chicken and poultry products, provides a vehicle for global spread. In this  
73 case, local phylogeographic signals can be weakened not only by the rapid movement of  
74 lineages around the world, but also by genomic changes that occur within the reservoir host.  
75 This may make it difficult to attribute the country of origin based on the *Campylobacter*  
76 isolate genome alone. Sequence based analyses have shown that populations of the main  
77 human disease-causing *Campylobacter* species, *C. jejuni* and *C. coli*, are highly structured  
78 into clusters of related lineages, known as clonal complexes, that share four or more alleles at  
79 7 multi-locus sequence typing (MLST) level (Dingle et al., 2005, Sheppard et al., 2010b). In  
80 *C. jejuni*, host-associated clonal complexes can be identified based upon the frequency with

81 which particular genotypes are isolated from different hosts (Sheppard et al., 2011, Sheppard  
82 et al., 2014). Many of these lineages are globally distributed (Sheppard et al., 2010a) but  
83 despite this strong host signal, there is evidence for phylogeographic structuring and the  
84 proliferation of distinct lineages in different countries (McTavish et al., 2008, Asakura et al.,  
85 2012).

86

87 Horizontal gene transfer in recombining bacteria, such as *Campylobacter* (Sheppard et al.,  
88 2008, Wilson et al., 2008, Sheppard et al., 2013a), can provide information about ecological  
89 differences between lineages. For example, when a lineage transfers to a new animal host it  
90 may acquire DNA from the resident population by HGT. This has been shown in host  
91 generalist *Campylobacter jejuni* lineages isolated from chicken that sometimes contain alleles  
92 that originated in chicken-specialist genotypes (McCarthy et al., 2007, Wilson et al., 2008).  
93 Here we applied comparable approaches to investigate if HGT can lead to signatures of  
94 recombination that discriminate isolates from North America and the UK using genomic data.  
95 Using matched pairs of Canadian and UK isolates, we identify genes that are prone to  
96 recombination, and will therefore pick up a local DNA more rapidly, and hypothesise that  
97 these genes may acquire a biogeographical signal.

98

99 **Materials and Methods**

100

101 ***Bacterial Isolates and Genome Sequencing***

102 A total of 294 sequenced isolates were analysed, of which 131 genomes were generated in  
103 this study, augmented by 163 previously published genomes (Sheppard et al., 2014, Sheppard  
104 et al., 2013a, Sheppard et al., 2013b). Sequencing reads for all genomes studied are available  
105 from the NCBI short read archive associated with BioProject: PRJNA312235 (Individual  
106 SRA accession numbers can be found in Table S1).

107

108 **Canadian isolates:** Isolates were collected from chicken and bovine faecal samples between  
109 **July 2004** and July 2006 from farms at diverse locations in **Alberta**. Samples were placed on  
110 ice and processed within 6 h as described by (Jokinen et al., 2010). Approximately 5 g of  
111 faecal matter was mixed with 5 mL of phosphate buffered saline (PBS) to form uniform  
112 slurry. One-millilitre aliquots of the PBS-faecal samples were added to 20 mL of Bolton  
113 broth containing 5% (v/v) lysed horse blood and selective supplement (Diergaardt et al.,  
114 2004) and incubated at 42°C for 24 h under microaerophilic conditions prior to plating 20 µl  
115 onto supplemented charcoal cefoperazone deoxycholate agar (CCDA). The plates were  
116 incubated for a further 48h at 42°C. Human samples were acquired from clinical laboratories  
117 in three Canadian provinces. These were re-plated from frozen glycerol stocks and the DNA  
118 extracted as described below.

119

120 Presumptive *Campylobacter* colonies were cultured onto blood agar plates and tested using  
121 biochemical oxidase and catalase tests. A multiplex PCR assay was used to detect 16S rRNA  
122 gene sequences and *C. jejuni* and *C. coli* specific primers *mapA* and *ceuE*, respectively

**Commented [BP1]:** Emma, Ed & Cathy to check details.

**Commented [ELS2]:** Yes, all chicken and cattle samples were collected in Alberta and dates entered.

123 (Denis et al., 1999). Positive *Campylobacter* isolates were sub-cultured on Mueller-Hinton  
124 agar and genomic DNA was extracted using the Wizard Genomic DNA Purification Kit as  
125 per manufacturer's instructions (Promega). DNA integrity was checked on an agarose gel and  
126 purity and concentration determined by optical density. Purified genomic DNA was sent to  
127 Canada's Michael Smith Genome Sciences Centre (Vancouver, Canada) and sequenced using  
128 the Illumina HiSeq 2000 platform.

129

130 **American isolates:** Isolates were collected from cattle faecal samples between December  
131 2008 and June 2010 from diverse locations within the Salinas Valley watershed, California.  
132 Samples were placed on ice and processed within 12 h. Cattle faeces were inoculated into a  
133 six-well microtiter plate containing 6 ml 1× Anaerobe Basal Broth (Oxoid) amended with  
134 Preston supplement (when reconstituted consists of: amphotericin B (0.01 mg/ml), rifampicin  
135 (0.01mg/ml), trimethoprim lactate (0.01mg/ml), and polymixin B (5UI/ml) (Oxoid), using a  
136 sterile cotton swab. Microtiter plates were placed inside plastic ZipLoc bags filled with a  
137 microaerobic gas mixture (1.5% O<sub>2</sub>, 10% H<sub>2</sub>, 10% CO<sub>2</sub>, and 78.5% N<sub>2</sub>) and incubated for  
138 24 h at 37°C, while shaking at 40 rpm. Subsequently, 10-µl of these enrichment cultures were  
139 plated onto anaerobe basal agar (ABA) plates, amended with 5% laked horse blood and CAT  
140 supplement (cefoperazone (0.008mg/ml), amphotericin B (0.01 mg/ml), and teicoplanin  
141 (0.004 mg/ml) (Oxoid)). All plates were then incubated under microaerobic conditions at  
142 37°C for 24 h. Bacterial cultures were passed through 0.2 µm mixed cellulose ester filters  
143 onto ABA plates and incubated at 37°C under microaerobic conditions. After 24 h, single  
144 colonies were streaked onto fresh ABA plates and incubated 24–48 h at 37°C for purification.

145

**Commented [BP3]:** Craig et al to check details.



146 DNA was extracted from a pure culture colony using Wizard Genomic DNA Purification Kit  
147 (Promega, Madison, WI). *Campylobacter* species was designated by 16S rDNA sequencing,  
148 using the primer pairs as described by Lane (1991). Genome sequencing was performed on an  
149 Illumina MiSeq sequencer using the KAPA Low-Throughput Library Preparation Kit with  
150 Standard PCR Amplification Module (Kapa Biosystems, Wilmington, MA), following  
151 manufacturer's instructions except for the following changes; 750 ng DNA was sheared at 30  
152 psi for 40 s and size selected to 700–770 bp following Illumina protocols. Standard desalted  
153 TruSeq LT and PCR Primers were ordered from Integrated DNA Technologies (Coralville,  
154 IA) and used at 0.375 and 0.5  $\mu$ M final concentrations, respectively. PCR was reduced to 3–5  
155 cycles. Libraries were quantified using the KAPA Library Quantification Kit (Kapa), except  
156 with 10  $\mu$ l volume and 90-s annealing/extension PCR, then pooled and normalized to 4 nM.  
157 Pooled libraries were re-quantified by ddPCR on a QX200 system (Bio-Rad), using the  
158 Illumina TruSeq ddPCR Library Quantification Kit and following manufacturer's protocols,  
159 except with an extended 2-min annealing/extension time. Libraries were sequenced using 2  $\times$   
160 250 bp paired end v2 reagent kit on a MiSeq instrument (Illumina) at 13.5 pM, following  
161 manufacturer's protocols. Reads were obtained from SeqWright (Houston, TX).

162

163 **UK isolates:** Sequenced isolates from Canada and the USA were augmented by 163  
164 previously published *Campylobacter* genomes collected between 1980 and 2012 from a range  
165 of sources, including cattle (54), chicken (80), pig (9), environmental (49), wild bird species  
166 (12) and human clinical cases (73) (Sheppard et al., 2014, Sheppard et al., 2013a, Sheppard et  
167 al., 2013b).

168

169 **UK clinical test isolates:** In addition, 383 clinical samples collected from the John Radcliffe  
170 Hospital in Oxford between June and October 2011 were used as a test dataset to attribute  
171 source according to geography (Cody et al., 2013). These genomes were downloaded from  
172 <http://pubmlst.org/campylobacter/>.

173

#### 174 *Population structure and selection of isolate pairs*

175 Isolate genomes were archived in the web-accessible BIGSdb database that supports  
176 functionality for identifying gene presence and allelic variation, by comparison to a reference  
177 locus list (Jolley and Maiden, 2010, Sheppard et al., 2012, Meric et al., 2014). This list  
178 comprised 1,623 locus designations from the annotated genome of *C. jejuni* strain  
179 NCTC11168 (Genbank accession number: NC\_002163.1) (Gundogdu et al., 2007, Parkhill et  
180 al., 2000). Reference loci were identified in each of the 294 isolate genomes using BLAST.  
181 Loci were recorded as present if the sequence had  $\geq 70\%$  nucleotide identity over  $\geq 50\%$  of the  
182 gene length. Each gene was aligned individually using MAFFT (Katoh et al., 2002), and  
183 concatenated into a single multi-FASTA alignment file for each isolate for a total alignment  
184 of 1,585,605 bp. Phylogenetic trees were constructed from a whole-genome alignment of  
185 *C. jejuni* (n=229) and *C. coli* (n=55) isolates based on 103,878 and 806,657 variable sites,  
186 respectively, using an approximation of the maximum likelihood algorithm (Tamura et al.,  
187 2013, Kumar et al., 2016). UK isolates from matching hosts were paired with their closest  
188 match from Canada. In total, 15 pairs of isolates were matched by source host and clonal  
189 complex (Figure 1). All paired isolates shared 1,378 genes giving rise to a core-genome  
190 alignment of 1,287,560 bp.

191

#### 192 *Analysis of co-ancestry and inference of recombination hot regions*

**Commented [NM4]:** Maybe put the criteria in here (< 1200 bp differences across the 1378 genes) rather than results.

193 The predicted co-ancestry of the paired isolates was determined based on whole genome  
194 sequences using *fineSTRUCTURE* (Yahara et al., 2013) and visualized as a heat map (Figure  
195 2). This algorithm infers the number of clusters (K) and partitions the strains into K  
196 subgroups with indistinguishable genetic ancestry, based on the likelihood of co-ancestry  
197 using a Bayesian MCMC (Markov chain Monte Carlo) approach (Lawson et al., 2012).  
198 Previous estimates of recombination rate and generation time (Webb and Blaser, 2002,  
199 Wilson et al., 2009, Morelli et al., 2010) were used to prepare a recombination map file  
200 specifying the same recombination rate per-site per-generation of SNPs. The predicted  
201 ancestry of co-inherited SNPs or chunks was calculated and SNPs with uncertain estimates of  
202 their donor of more than 20 kb results were removed. The results were visualised in the  
203 UCSC (The University of California Santa Cruz) browser (Kent et al., 2002) (Figure 2).  
204 Following burn-in, Markov chain Monte Carlo (MCMC) iterations were run 100,000 times  
205 in *fineSTRUCTURE* (version 0.02) (Lawson et al., 2012) with a thinning interval of 100.  
206 Population assignments runs were performed twice.

207 ||  
208 The time to the most recent common ancestor (TMRCA) in each pair was estimated using the  
209 model described in Didelot et al., (2013) and summarised here briefly. Pairs of genomes  
210 share a common ancestor  $t$  years ago and have been subject to mutation at a rate  $\mu$  and  
211 recombination at rate  $\rho$ . The mutation rate of  $2.9 \times 10^{-5}$  per site per year was used as reported  
212 in Sheppard et al., (2010b), which is similar to the rates estimated in Wilson et al., (2008,  
213 2009) ). The effect of recombination is to introduce a high density of polymorphism similar  
214 to the ClonalFrame model (cite) but with the advantage that this density can vary between  
215 recombination events to reflect differences in evolutionary distance between donors and

**Commented [XD5]:** I think this whole paragraph was copied by mistake and needs to be removed since here we do not use simulations

**Commented [XD6]:** Didelot X, Falush D (2007) Inference of bacterial microevolution using multilocus sequence data. *Genetics* 175:1251–66. doi: 10.1534/genetics.106.063305

Didelot X, Wilson DJ (2015) ClonalFrameML: Efficient Inference of Recombination in Whole Bacterial Genomes. *PLOS Comput Biol* 11:e1004041. doi: 10.1371/journal.pcbi.1004041

216 recipients(cite). In each pairwise comparison, the TMRCA and recombination rate parameters  
217 are estimated based on a core genome alignment, with 95% credibility intervals (Table 2).

218

### 219 *Epidemiological markers of geographical clustering*

220 Neighbour-joining phylogenetic trees were constructed for all genes that demonstrated  
221 pairwise diversity above 2% nucleotide diversity (Table S2). Individual gene phylogenies  
222 were constructed in MEGA for all 57 genes. Isolates were assigned to a putative source  
223 population based on the seven highly recombining genes that showed the greatest level of  
224 clustering by geography. Probabilistic assignment of geographical source is based on the  
225 allele frequencies in the reference population data sets for each of the seven loci. This  
226 analysis was performed using Structure, a Bayesian model-based clustering method designed  
227 to infer population structure and assign individuals to populations using multilocus genotype  
228 data (Sheppard et al., 2010a, Pritchard et al., 2000). Canadian and USA isolates were  
229 combined as a North American population for comparison with UK isolates (Table S1).

230

### 231 *Attribution of clinical isolates to country based on 7 geographically segregating genes*

232 The source attribution model was tested with isolates of a known source. Self-assignment of a  
233 random subset of the comparison dataset was conducted by removing a third of the isolates  
234 from each candidate population (n = 73). The remainder were used as the reference set (78  
235 North American isolates to compare with 68 UK isolates). Structure was run for 100,000  
236 iterations following a burn-in period of 10,000 iterations using the no admixture model to  
237 assign individuals to putative populations. The assignment probability for each source was  
238 calculated for each isolate individually and were attributed to origin populations when the  
239 attribution probability was greater than 0.50.

240

**Commented [XD7]:** Morelli G, Didelot X, Kusecek B, et al (2010) Microevolution of *Helicobacter pylori* during prolonged infection of single hosts and within families. PLoS Genet 6:e1001036. doi: 10.1371/journal.pgen.1001036

Didelot X, Nell S, Yang I, et al (2013) Genomic evolution and transmission of *Helicobacter pylori* in two South African families. Proc Natl Acad Sci USA 110:13880–13885. doi: 10.1073/pnas.1304681110

241 **Results**

242 Core genomes of isolates from North America and the UK were compared, and there was no  
243 observable clustering by country or continent on a neighbour-joining tree (Figure 1). STs  
244 sampled in both *Campylobacter* populations belonged to clonal complexes that can be  
245 classified as specialist and host generalist based upon the frequency at which they have been  
246 isolated from different hosts. These included chicken specialist sequence types ST-257, ST-  
247 283, ST-353, ST-354, ST-443, ST-573, ST-574 and ST-661 clonal complexes, cattle  
248 specialist ST-61 and ST-42 clonal complexes, and host generalist ST-21, ST-45, ST-206 and  
249 ST-48 complexes (Figure 1 and Table S1).

250

251 ***Matched isolates share more common ancestry with isolates from the same country***

252 To minimise the effect of host adaptation and maximize the opportunity of identifying  
253 genetic signatures of geographic separation, a subset of 15 isolate pairs were chosen based  
254 upon their phylogenetic clustering, < 1,200 bp difference in 1,378 core genome loci. In each  
255 case, isolate pairs contained one Canadian and one UK isolate of the same clonal complex  
256 sampled from the same host species (Table 1). The predicted co-ancestry of the paired  
257 isolates was calculated based on core genome alignments using *fineSTRUCTURE* (Yahara et  
258 al., 2013) (Figure 2). DNA sequence haplotype regions were coloured by predicted  
259 inheritance from donor isolates and the average frequency of co-ancestry of DNA ‘chunks’  
260 from isolates within the same country (0.58) was significantly greater than that for isolates  
261 from different countries (0.32). The degree of inheritance for each gene was calculated and  
262 genes that have been predicted to inherit the most DNA from donor isolates of different  
263 countries was surmised (Table S2).

264

Commented [B8]: Koji - Pvalue and test used

265 ***Matched isolates share recent common ancestors but have since experienced significant***  
266 ***recombination***

267 The estimated time since the most recent common ancestor (TMRCA) was calculated for  
268 each UK/American pair of genomes as previously described (Didelot et al., 2013), using the  
269 mutation rate of  $2.9 \times 10^{-5}$  per site per year reported in Sheppard et al. (2010b), which is  
270 consistent with estimates in Wilson et al. (2008, 2009). In each pairwise comparison, the  
271 level of divergence along the genome (Figure 3) was used to estimate the TMRCA and  
272 recombination rate, with 95% credibility intervals around these parameters (Table 2). All  
273 pairs were estimated to have shared ancestors between one and five years ago, with two  
274 exceptions, namely the two *C. coli* pairs, for which the TMRCA was around 25 years ago.

275 The ratio  $r/m$  of rates at which recombination and mutation introduce polymorphism (cite)  
276 was estimated to be around 20-30 except in the two *C. coli* pairs with larger TMRCA, for  
277 which a much smaller value was estimated around  $r/m=4$ .

279 ***Highly recombining genes as markers of geographical attribution***

280 A pairwise comparison of the matched pairs was used to quantify the level of divergence in  
281 each gene within the core genome (1,147 genes) of the paired isolates. Most genes showed  
282 low diversity, indicative of closely related pairs. Polymorphism in genes with less than 2%  
283 divergence between pairs (white and red in Figure 3) are likely to be the result of mutation or  
284 recombination with a tract of DNA with high nucleotide identity, so that only one or two  
285 substitutions are visible. Genes with greater than 2% divergence between pairs are likely to  
286 have recombined as numerous substitutions have been introduced (blue in Figure 3). Fifty-  
287 seven genes (e.g. *Cj0034c* and *Cj0635*) had a high level (>2%) of nucleotide divergence and  
288 high probability of recombination in all 15 pairs. This result did not arise just by chance:

**Commented [NM9]:** The approx. tenfold longer TMRCA of *C. coli* vs *C. jejuni* pairs is so striking – and within pairs based on the 1200bp difference cut off that may operate differently across species, and a much less densely sampled *C. coli* population to generate sampling of close isolates. So many reasons may be contributing to this apparent difference even if true. To me these questions treating these *coli* and *jejuni* pairs as dealing with the same thing.

**Commented [XD10]:** Vos M, Didelot X (2009) A comparison of homologous recombination rates in bacteria and archaea. ISME J 3:199–208. doi: 10.1038/ismej.2008.93

289 overall recombination was inferred in around 25% of the genes in each pair and so if  
290 recombination was random, the probability that all 15 pairs had recombined for a given gene  
291 would be extremely small ( $0.25^{15}=9.3 \times 10^{-10}$ ).

292

293 Individual gene trees were generated for these 57 genes from which the most recombination  
294 could be identified. The seven genes that gave the clearest geographic clustering were used  
295 for further analysis of geographical attribution using Structure as previously described  
296 (Sheppard et al., 2010a, Pritchard et al., 2000). A self-test was performed on our collection of  
297 xxx isolates and in 76.7% of cases the source continent was correctly attributed. The  
298 percentages of correctly attributed isolates by population were not significantly different, at  
299 76.9% for North America and 76.5% for the UK. Where an isolate was incorrectly attributed  
300 to a population there was a higher average reported attribution probability (0.85) in the case  
301 of UK isolates compared with North American isolates (0.67). The proportion of UK isolates  
302 correctly attributed to the UK reference population was 70%, while the proportion of North  
303 American isolates correctly attributed was slightly higher at 76%.

304

#### 305 *Attribution of clinical isolates to country based on seven selected genes*

306 The same geographical attribution model was applied to 383 clinical *C. jejuni* isolates from  
307 the Oxfordshire *Campylobacter* Surveillance Study in the UK, accessed via  
308 [pubMLST.org/campylobacter](http://pubMLST.org/campylobacter), and for which details of recent foreign travel were provided  
309 (Cody et al., 2013). The model correctly assigned 34 of the 46 (73.9%) isolates where recent  
310 foreign travel had previously been declared, to a non-UK source of origin (Figure 4). In total,  
311 approximately half (47%) of the collected clinical isolates could be attributed to the UK.

312

313

**Commented [ET11]:** Ben, did you ever have a chance to compare the performance of these 7 genes vs. the geographical attribution estimates obtained with genes with lesser geographical signal? It might be nice to have an idea of how much improvement in performance was obtained.

**Commented [NM12]:** This looks very good but looking at the figure I think that it is true to say that about 50% of those without a travel history would also have been assigned as non-UK. I think that it would be informative to give the figure for those without a travel history.

314 **Discussion**

315 In gut dwelling bacteria, isolation in different host species, and barriers to recombination  
316 between populations, overtime, can lead to population differentiation reflected in the genome.

317 In *C. jejuni*, this can be seen at different levels; the proliferation of certain lineages to a  
318 particular host species, that are abundant in one host and rare or absent in others (Sheppard et  
319 al., 2011, Griekspoor et al., 2013, Sheppard et al., 2010a), secondly, as the increased  
320 frequency of host associated nucleotide substitutions in multiple lineages (that reflect  
321 adaptation to the host) drift in physically isolated populations (Sheppard et al., 2013b). This  
322 host-associated genetic structuring can be informative for understanding the evolution of *C.*  
323 *jejuni* (Dearlove et al., 2016), but can also be used in a more practical way to identify the  
324 source of isolates causing human infection by identifying genomic signatures (resulting from  
325 adaptation or drift) in the infecting isolate that are associated with populations in particular  
326 reservoir hosts (Sheppard et al., 2009, Wilson et al., 2008). Quantitative source attribution  
327 models, based upon the probability that a particular clinical isolate originated in different  
328 reservoirs, has been widely used to estimate the risk of human infection from different food  
329 production animals and other sources (Colles et al., 2008, French et al., 2005, Mullner et al.,  
330 2009, Sheppard et al., 2009, Roux et al., 2013, Griekspoor et al., 2013, Viswanathan et al.,  
331 2016) and have informed intervention strategies and public health policy (Cody et al., 2013,  
332 Cody et al., 2012).

333

334 The accuracy of probabilistic source attribution models is influenced by the degree of which  
335 indicative markers in the isolate genome, such as MLST locus alleles, can be placed within a  
336 source population. This would be relatively straightforward for markers that segregate  
337 absolutely by source, but in *C. jejuni* and *C. coli* it is common that alleles are present in more

**Commented [NM13]:** Nothing particular to gut dwelling bacteria as opposed to e.g. maxillary sinus dwelling bacteria etc.



338 than one population, but at different frequencies. In simple attribution models using MLST  
339 data, *C. jejuni* and *C. coli* isolates from chickens in the Netherlands, Senegal and the USA  
340 have been more closely related to UK chicken isolate populations rather than to populations  
341 from other host species in the same country (Sheppard et al., 2010a). While genomic  
342 signatures of host association can transcend geographic structuring within *C. jejuni* and *C.*  
343 *coli* populations, there can be differences in the genotypes that are isolated from different  
344 countries (Mohan et al., 2013, Asakura et al., 2012, Kivisto et al., 2014, Islam et al., 2014,  
345 Prachantasena et al., 2016). This presents challenges, not only for attributing the source of  
346 infections among travellers returning from foreign locations (Mughini-Gras et al., 2014), but  
347 also for understanding disease epidemiology in the context of a global food industry.

348

349 Following the occupation of a new niche *C. jejuni* and *C. coli* can acquire DNA signatures  
350 through recombination (Wilson et al., 2009, Sheppard et al., 2013a, Sheppard et al., 2008)  
351 and local DNA signatures via HGT, from resident strains. To quantify the extent to which  
352 isolates from the same country share DNA sequence, we compared 15 isolate pairs from  
353 different countries, that to minimise the effect of clonal inheritance and host-associated  
354 variation were matched by both clonal complex and source -. The predicted ancestry of co-  
355 inherited SNPs was nearly twice as high among isolates from same country compared to  
356 those from different countries. While this represents a relatively weak signal of geographic  
357 association, compared to host association, there was a quantifiable local (national) signal that  
358 can be used to investigate geographical clustering.

359

360 Since, recombination introduces more nucleotide substitutions than during mutation in *C.*  
361 *jejuni* and *C. coli* (Webb and Blaser, 2002, Wilson et al., 2009, Morelli et al., 2010), genes

362 with evidence of elevated recombination rates, that share a gene pool, will more rapidly  
363 acquire local signals of sequence variation than genes with lower recombination rates. These  
364 genes represent potential targets for use as biogeographic epidemiological markers. Pairwise  
365 isolate comparison revealed that nucleotide divergence was <2% across the majority of the  
366 genome (Figure 2, Table S2). However, some genes consistently had more sequence variation  
367 in multiple isolate pairs, potentially indicating enhanced recombination at these loci.

368

369 Several of these genes have been annotated with functions associated with DNA processing,  
370 transcription, repair and maintenance. This may reflect the mechanisms of recombination and  
371 horizontal gene transfer. Other genes with evidence of elevated recombination included those  
372 associated with surface exposed proteins with roles in glycosylation, motility and secretion  
373 which would form part of an initial interaction with the host/environment (Table S2). The *C.*  
374 *jejuni* N-acetyltransferase PseH (Cj1313) plays a key role in O-linked glycosylation, which  
375 contributes to flagellar formation, motility and pseudoamino acid biosynthesis (Song et al.,  
376 2015, McNally et al., 2006) and is important in host colonisation (Guerry et al., 2006). The  
377 variable outer membrane protein gene *PorA*, which has been used as part of extended MLST  
378 schemes (Dingle et al., 2008, Cody et al., 2009) was also among those genes with evidence of  
379 elevated recombination. This may explain why weak allopatric signals have been associated  
380 with sequence variation in the *PorA* gene in addition to source attribution signals (Sheppard  
381 et al., 2010a, Smid et al., 2013, Mughini-Gras et al., 2014).

382

383 Three efflux pump genes *Cj0034c*, *Cj0619* and *Cj1174* genes, that have been implicated in  
384 fluoroquinolone resistance, showed elevated recombination and phylogeographic variation  
385 (Table S2) (Luangtongkum et al., 2009, Ge et al., 2005). Clinical and agricultural prescription

**Commented [NM14]:** Here you cover some genes used for geography and others not. Good to make explicit. Only the supplementary material shows the 7 genes selected and this is reasonably central and would be better included even as a sentence in the paper. More generally I wonder about you considering the criteria for selecting genes along the lines of (1) first principles that substantial recombination needed to allow generation of a signal. (2) some biological processes e.g. Cipro resistance likely to support geographical patterning and others e.g. "initial interaction with host" not likely to be geographically informative etc. but could appear to be if e.g. more samples from a particular place were from a particular species, and that your selection was based on empirical association.

386 of broad-spectrum antibiotics such as quinolones varies worldwide. Since the late 1990's the  
387 agricultural use of fluoroquinolones has declined following governmental intervention in  
388 Europe and North America (Chang et al., 2015, Nelson et al., 2007). However, resistant  
389 isolates remain common and the level of resistance can vary from country to country (Pham  
390 et al., 2015). Higher levels of fluoroquinolone resistance has been observed among isolates  
391 from patients returning from foreign travel (Gaudreau et al., 2014). The identification of  
392 genes associated with efflux pumps, among those with high levels of inferred recombination,  
393 suggests a role in the emergence of fluoroquinolones-resistance and provides a useful  
394 indicator for geographic segregation of isolates.

395  
396 Using signatures of local recombination in *Campylobacter* genomes, has the potential to  
397 identify the country of origin and attribute the source of infection, among returning  
398 travellers. In this study 74% of isolates from individuals, that had declared recent foreign  
399 travel, were attributed to non-UK sources. However, in the absence of genetic elements that  
400 segregate absolutely by geography, the model relies upon the availability of large reference  
401 datasets from reservoir populations in different countries for frequency-dependent attribution.  
402 Although this limits the applicability of the approach using currently available data the  
403 statistical genetics methodologies employed here provide a quantitative means for identifying  
404 genomic signatures of allopatry. This potentially enables the evaluation of transmission  
405 dynamics through global livestock trade networks. *Campylobacter* populations are highly  
406 structured with some lineages having greater significance in human disease than others, either  
407 because of enhanced capacity to survive through slaughter and food production [Ref] or  
408 increased antimicrobial resistance (Wimalaratna et al., 2013, Cody et al., 2010). Monitoring

Commented [B15]: Yahara and Meric when published (In press, EM)

409 the spread of these strains may be useful for evidence-based interventions targeting strains  
410 that are a significant global health burden.

411

412 **Acknowledgements:**

413 SKS is a WellcomeTrust Fellow (088786/C/09/Z) and research in his laboratory is funded by  
414 grants from the Medical Research Council (MR/L015080/1), the Food Standards Agency  
415 (FS246004) and the Biotechnology and Biological Sciences Research Council  
416 (BB/I02464X/1). We acknowledge Canada's Michael Smith Genome Sciences Centre,  
417 Vancouver, Canada and USA sequencing for sequencing. GM is supported by a postdoctoral  
418 research fellowship from Health and care research, Wales (HF-14-13). KY was supported by  
419 a JSPS Research Fellowships for Young Scientists. We also acknowledge MRC CLIMB and  
420 HPC Wales for the use of HPC facilities.

421

422 **Data Accessibility**

423 Draft assembly genomes and short sequencing reads generated in this study have been  
424 deposited in NCBI GenBank database and/or the Short Read Archive associated with  
425 BioProject(s): PRJNA312235. Individual accession numbers can be found in table S1.

426

427 **Author contributions**

428 BP, GM, XD and SKS designed research; BP, GM, KY, HW, SM and XD performed  
429 research; BP, GM, KY, HW, SM, NM, XD, CTP and SKS analysed results; MDH, ELS,  
430 CDC, ENT, KKC, SH, WGM, AJC, KAJ, MCJM, NM and SKS provided isolates, genomes  
431 or software and BP, GM, CTP and SKS wrote the manuscript.

432

433 **Conflict of Interest Statement**

434 The authors declare no competing interests.

435

436 **References**

- 437 ACHTMAN, M. 2008. Evolution, population structure, and phylogeography of genetically  
438 monomorphic bacterial pathogens. *Annu Rev Microbiol*, 62, 53-70.
- 439 ASAKURA, H., BRUGGEMANN, H., SHEPPARD, S. K., EKAWA, T., MEYER, T. F., YAMAMOTO, S. & IGIMI,  
440 S. 2012. Molecular evidence for the thriving of *Campylobacter jejuni* ST-4526 in Japan. *PLoS*  
441 *One*, 7, e48394.
- 442 CHANG, Q., WANG, W., REGEV-YOCHAY, G., LIPSITCH, M. & HANAGE, W. P. 2015. Antibiotics in  
443 agriculture and the risk to human health: how worried should we be? *Evol Appl*, 8, 240-7.
- 444 CODY, A. J., CLARKE, L., BOWLER, I. C. & DINGLE, K. E. 2010. Ciprofloxacin-resistant  
445 campylobacteriosis in the UK. *Lancet*, 376, 1987.
- 446 CODY, A. J., MAIDEN, M. J. & DINGLE, K. E. 2009. Genetic diversity and stability of the *porA* allele as a  
447 genetic marker in human *Campylobacter* infection. *Microbiology*, 155, 4145-54.
- 448 CODY, A. J., MCCARTHY, N. D., JANSEN VAN RENSBURG, M., ISINKAYE, T., BENTLEY, S. D., PARKHILL,  
449 J., DINGLE, K. E., BOWLER, I. C., JOLLEY, K. A. & MAIDEN, M. C. 2013. Real-time genomic  
450 epidemiological evaluation of human *Campylobacter* isolates by use of whole-genome  
451 multilocus sequence typing. *J Clin Microbiol*, 51, 2526-34.
- 452 CODY, A. J., MCCARTHY, N. M., WIMALARATHNA, H. L., COLLES, F. M., CLARK, L., BOWLER, I. C.,  
453 MAIDEN, M. C. & DINGLE, K. E. 2012. A longitudinal 6-year study of the molecular  
454 epidemiology of clinical campylobacter isolates in Oxfordshire, United kingdom. *J Clin*  
455 *Microbiol*, 50, 3193-201.
- 456 COLLES, F. M., JONES, T. A., MCCARTHY, N. D., SHEPPARD, S. K., CODY, A. J., DINGLE, K. E., DAWKINS,  
457 M. S. & MAIDEN, M. C. 2008. Campylobacter infection of broiler chickens in a free-range  
458 environment. *Environ Microbiol*, 10, 2042-50.
- 459 COMAS, I., COSCOLLA, M., LUO, T., BORRELL, S., HOLT, K. E., KATO-MAEDA, M., PARKHILL, J., MALLA,  
460 B., BERG, S. & THWAITES, G. 2013. Out-of-Africa migration and Neolithic coexpansion of  
461 *Mycobacterium tuberculosis* with modern humans. *Nature genetics*, 45, 1176-1182.
- 462 DEARLOVE, B. L., CODY, A. J., PASCOE, B., MERIC, G., WILSON, D. J. & SHEPPARD, S. K. 2016. Rapid  
463 host switching in generalist *Campylobacter* strains erodes the signal for tracing human  
464 infections. *ISME J*, 10, 721-9.
- 465 DENIS, M., SOUMET, C., RIVOAL, K., ERMEL, G., BLIVET, D., SALVAT, G. & COLIN, P. 1999.  
466 Development of a m-PCR assay for simultaneous identification of *Campylobacter jejuni* and  
467 *C. coli*. *Lett Appl Microbiol*, 29, 406-10.
- 468 DIDELOT, X., LAWSON, D., DARLING, A. & FALUSH, D. 2010. Inference of homologous recombination  
469 in bacteria using whole-genome sequences. *Genetics*, 186, 1435-49.
- 470 DIDELOT, X., NELL, S., YANG, I., WOLTEMATE, S., VAN DER MERWE, S. & SUERBAUM, S. 2013.  
471 Genomic evolution and transmission of *Helicobacter pylori* in two South African families.  
472 *Proc Natl Acad Sci U S A*, 110, 13880-5.
- 473 DIERGAARDT, S. M., VENTER, S. N., SPREETH, A., THERON, J. & BROZEL, V. S. 2004. The occurrence of  
474 campylobacters in water sources in South Africa. *Water Res*, 38, 2589-95.
- 475 DINGLE, K. E., COLLES, F. M., FALUSH, D. & MAIDEN, M. C. 2005. Sequence typing and comparison of  
476 population biology of *Campylobacter coli* and *Campylobacter jejuni*. *J Clin Microbiol*, 43, 340-  
477 7.
- 478 DINGLE, K. E., MCCARTHY, N. D., CODY, A. J., PETO, T. E. & MAIDEN, M. C. 2008. Extended sequence  
479 typing of *Campylobacter* spp., United Kingdom. *Emerg Infect Dis*, 14, 1620-2.
- 480 FALUSH, D., WIRTH, T., LINZ, B., PRITCHARD, J. K., STEPHENS, M., KIDD, M., BLASER, M. J., GRAHAM,  
481 D. Y., VACHER, S., PEREZ-PEREZ, G. I., YAMAOKA, Y., MEGRAUD, F., OTTO, K., REICHARD, U.,  
482 KATZOWITSCH, E., WANG, X., ACHTMAN, M. & SUERBAUM, S. 2003. Traces of human  
483 migrations in *Helicobacter pylori* populations. *Science*, 299, 1582-5.

484 FRENCH, N., BARRIGAS, M., BROWN, P., RIBIERO, P., WILLIAMS, N., LEATHERBARROW, H., BIRTLES,  
485 R., BOLTON, E., FEARNHEAD, P. & FOX, A. 2005. Spatial epidemiology and natural population  
486 structure of *Campylobacter jejuni* colonizing a farmland ecosystem. *Environ Microbiol*, 7,  
487 1116-26.

488 GAGNEUX, S. & SMALL, P. M. 2007. Global phylogeography of *Mycobacterium tuberculosis* and  
489 implications for tuberculosis product development. *Lancet Infect Dis*, 7, 328-37.

490 GAUDREAU, C., BOUCHER, F., GILBERT, H. & BEKAL, S. 2014. Antimicrobial susceptibility of  
491 *Campylobacter jejuni* and *Campylobacter coli* isolates obtained in Montreal, Quebec,  
492 Canada, from 2002 to 2013. *J Clin Microbiol*, 52, 2644-6.

493 GE, B., MCDERMOTT, P. F., WHITE, D. G. & MENG, J. 2005. Role of efflux pumps and topoisomerase  
494 mutations in fluoroquinolone resistance in *Campylobacter jejuni* and *Campylobacter coli*.  
495 *Antimicrob Agents Chemother*, 49, 3347-54.

496 GRIEKSPoor, P., COLLES, F. M., MCCARTHY, N. D., HANSBRO, P. M., ASHHURST-SMITH, C., OLSEN, B.,  
497 HASSELQUIST, D., MAIDEN, M. C. & WALDENSTROM, J. 2013. Marked host specificity and  
498 lack of phylogeographic population structure of *Campylobacter jejuni* in wild birds. *Mol Ecol*,  
499 22, 1463-72.

500 GUERRY, P., EWING, C. P., SCHIRM, M., LORENZO, M., KELLY, J., PATTARINI, D., MAJAM, G.,  
501 THIBAUT, P. & LOGAN, S. 2006. Changes in flagellin glycosylation affect *Campylobacter*  
502 autoagglutination and virulence. *Mol Microbiol*, 60, 299-311.

503 GUNDOGDU, O., BENTLEY, S. D., HOLDEN, M. T., PARKHILL, J., DORRELL, N. & WREN, B. W. 2007. Re-  
504 annotation and re-analysis of the *Campylobacter jejuni* NCTC11168 genome sequence. *BMC*  
505 *Genomics*, 8, 162.

506 ISLAM, Z., VAN BELKUM, A., WAGENAAR, J. A., CODY, A. J., DE BOER, A. G., SARKER, S. K., JACOBS, B.  
507 C., TALUKDER, K. A. & ENDTZ, H. P. 2014. Comparative population structure analysis of  
508 *Campylobacter jejuni* from human and poultry origin in Bangladesh. *Eur J Clin Microbiol*  
509 *Infect Dis*, 33, 2173-81.

510 JOKINEN, C. C., SCHREIER, H., MAURO, W., TABOADA, E., ISAAC-RENTON, J. L., TOPP, E., EDGE, T.,  
511 THOMAS, J. E. & GANNON, V. P. 2010. The occurrence and sources of *Campylobacter* spp.,  
512 *Salmonella enterica* and *Escherichia coli* O157:H7 in the Salmon River, British Columbia,  
513 Canada. *J Water Health*, 8, 374-86.

514 JOLLEY, K. A. & MAIDEN, M. C. 2010. BIGSdb: Scalable analysis of bacterial genome variation at the  
515 population level. *BMC Bioinformatics*, 11, 595.

516 KATOH, K., MISAWA, K., KUMA, K. & MIYATA, T. 2002. MAFFT: a novel method for rapid multiple  
517 sequence alignment based on fast Fourier transform. *Nucleic Acids Res*, 30, 3059-66.

518 KENT, W. J., SUGNET, C. W., FUREY, T. S., ROSKIN, K. M., PRINGLE, T. H., ZAHLER, A. M. & HAUSSLER,  
519 D. 2002. The human genome browser at UCSC. *Genome Res*, 12, 996-1006.

520 KIVISTO, R. I., KOVANEN, S., SKARP-DE HAAN, A., SCHOTT, T., RAHKIO, M., ROSSI, M. & HANNINEN,  
521 M. L. 2014. Evolution and comparative genomics of *Campylobacter jejuni* ST-677 clonal  
522 complex. *Genome Biol Evol*, 6, 2424-38.

523 KUMAR, S., STECHER, G. & TAMURA, K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis  
524 Version 7.0 for Bigger Datasets. *Mol Biol Evol*, 33, 1870-4.

525 LANE, D. J. 1991. 16S/23S rRNA sequencing. In: E., S. & M., G. (eds.) *Nucleic Acid Sequencing*  
526 *Techniques in Bacterial Systematics*. New York: Wiley and Sons.

527 LAWSON, D. J., HELLENTHAL, G., MYERS, S. & FALUSH, D. 2012. Inference of population structure  
528 using dense haplotype data. *PLoS Genet*, 8, e1002453.

529 LUANGTONGKUM, T., JEON, B., HAN, J., PLUMMER, P., LOGUE, C. M. & ZHANG, Q. 2009. Antibiotic  
530 resistance in *Campylobacter*: emergence, transmission and persistence. *Future Microbiol*, 4,  
531 189-200.

532 MCCARTHY, N. D., COLLES, F. M., DINGLE, K. E., BAGNALL, M. C., MANNING, G., MAIDEN, M. C. &  
533 FALUSH, D. 2007. Host-associated genetic import in *Campylobacter jejuni*. *Emerg Infect Dis*,  
534 13, 267-72.

535 MCNALLY, D. J., HUI, J. P., AUBRY, A. J., MUI, K. K., GUERRY, P., BRISSON, J. R., LOGAN, S. M. & SOO,  
536 E. C. 2006. Functional characterization of the flagellar glycosylation locus in *Campylobacter*  
537 *jejuni* 81-176 using a focused metabolomics approach. *J Biol Chem*, 281, 18489-98.

538 MCTAVISH, S. M., POPE, C. E., NICOL, C., SEXTON, K., FRENCH, N. & CARTER, P. E. 2008. Wide  
539 geographical distribution of internationally rare *Campylobacter* clones within New Zealand.  
540 *Epidemiol Infect*, 136, 1244-52.

541 MERIC, G., YAHARA, K., MAGEIROS, L., PASCOE, B., MAIDEN, M. C., JOLLEY, K. A. & SHEPPARD, S. K.  
542 2014. A reference pan-genome approach to comparative bacterial genomics: identification  
543 of novel epidemiological markers in pathogenic *campylobacter*. *PLoS One*, 9, e92798.

544 MOHAN, V., STEVENSON, M., MARSHALL, J., FEARNHEAD, P., HOLLAND, B. R., HOTTER, G. & FRENCH,  
545 N. P. 2013. *Campylobacter jejuni* colonization and population structure in urban populations  
546 of ducks and starlings in New Zealand. *Microbiologyopen*, 2, 659-73.

547 MOODLEY, Y., LINZ, B., YAMAOKA, Y., WINDSOR, H. M., BREUREC, S., WU, J. Y., MAADY, A.,  
548 BERNHOF, S., THIBERGE, J. M., PHUANUKOONNON, S., JOBB, G., SIBA, P., GRAHAM, D. Y.,  
549 MARSHALL, B. J. & ACHTMAN, M. 2009. The peopling of the Pacific from a bacterial  
550 perspective. *Science*, 323, 527-30.

551 MORELLI, G., DIDELOT, X., KUSECEK, B., SCHWARZ, S., BAHLOWANE, C., FALUSH, D., SUERBAUM, S. &  
552 ACHTMAN, M. 2010. Microevolution of *Helicobacter pylori* during prolonged infection of  
553 single hosts and within families. *PLoS Genet*, 6, e1001036.

554 MUGHINI-GRAS, L., SMID, J. H., WAGENAAR, J. A., A, D. E. B., HAVELAAR, A. H., FRIESEMA, I. H.,  
555 FRENCH, N. P., GRAZIANI, C., BUSANI, L. & VAN PELT, W. 2014. *Campylobacteriosis* in  
556 returning travellers and potential secondary transmission of exotic strains. *Epidemiol Infect*,  
557 142, 1277-88.

558 MULLNER, P., SPENCER, S. E., WILSON, D. J., JONES, G., NOBLE, A. D., MIDWINTER, A. C., COLLINS-  
559 EMERSON, J. M., CARTER, P., HATHAWAY, S. & FRENCH, N. P. 2009. Assigning the source of  
560 human *campylobacteriosis* in New Zealand: a comparative genetic and epidemiological  
561 approach. *Infect Genet Evol*, 9, 1311-9.

562 NELSON, J. M., CHILLER, T. M., POWERS, J. H. & ANGULO, F. J. 2007. Fluoroquinolone-resistant  
563 *Campylobacter* species and the withdrawal of fluoroquinolones from use in poultry: a public  
564 health success story. *Clin Infect Dis*, 44, 977-80.

565 PARKHILL, J., WREN, B. W., MUNGALL, K., KETLEY, J. M., CHURCHER, C., BASHAM, D.,  
566 CHILLINGWORTH, T., DAVIES, R. M., FELTWELL, T., HOLROYD, S., JAGELS, K., KARLYSHEV, A.  
567 V., MOULE, S., PALLAN, M. J., PENN, C. W., QUAIL, M. A., RAJANDREAM, M. A., RUTHERFORD,  
568 K. M., VAN VLIET, A. H., WHITEHEAD, S. & BARRELL, B. G. 2000. The genome sequence of the  
569 food-borne pathogen *Campylobacter jejuni* reveals hypervariable sequences. *Nature*, 403,  
570 665-8.

571 PHAM, N. T., THONGPRACHUM, A., TRAN, D. N., NISHIMURA, S., SHIMIZU-ONDA, Y., TRINH, Q. D.,  
572 KHAMRIN, P., UKARAPOL, N., KONGSRICHAROERN, T., KOMINE-AIZAWA, S., OKITSU, S.,  
573 MANEEKARN, N., HAYAKAWA, S. & USHIJIMA, H. 2015. Antibiotic Resistance of  
574 *Campylobacter jejuni* and *C. coli* Isolated from Children with Diarrhea in Thailand and Japan.  
575 *Jpn J Infect Dis*.

576 PRACHANTASENA, S., CHARUNUNTAKORN, P., MUANGNOICHAROEN, S., HANKLA, L., TECHAWAL, N.,  
577 CHAVEERACH, P., TUITEMWONG, P., CHOKESAJJAWATEE, N., WILLIAMS, N., HUMPHREY, T.  
578 & LUANGTONGKUM, T. 2016. Distribution and Genetic Profiles of *Campylobacter* in  
579 Commercial Broiler Production from Breeder to Slaughter in Thailand. *PLoS One*, 11,  
580 e0149585.



581 PRITCHARD, J. K., STEPHENS, M. & DONNELLY, P. 2000. Inference of population structure using  
582 multilocus genotype data. *Genetics*, 155, 945-59.

583 ROUX, F., SPROSTON, E., ROTARIU, O., MACRAE, M., SHEPPARD, S. K., BESSELL, P., SMITH-PALMER,  
584 A., COWDEN, J., MAIDEN, M. C., FORBES, K. J. & STRACHAN, N. J. 2013. Elucidating the  
585 aetiology of human *Campylobacter coli* infections. *PLoS One*, 8, e64504.

586 SHEPPARD, S. K., CHENG, L., MERIC, G., DE HAAN, C. P., LLARENA, A. K., MARTTINEN, P., VIDAL, A.,  
587 RIDLEY, A., CLIFTON-HADLEY, F., CONNOR, T. R., STRACHAN, N. J., FORBES, K., COLLES, F. M.,  
588 JOLLEY, K. A., BENTLEY, S. D., MAIDEN, M. C., HANNINEN, M. L., PARKHILL, J., HANAGE, W. P.  
589 & CORANDER, J. 2014. Cryptic ecology among host generalist *Campylobacter jejuni* in  
590 domestic animals. *Mol Ecol*, 23, 2442-51.

591 SHEPPARD, S. K., COLLES, F., RICHARDSON, J., CODY, A. J., ELSON, R., LAWSON, A., BRICK, G.,  
592 MELDRUM, R., LITTLE, C. L., OWEN, R. J., MAIDEN, M. C. & MCCARTHY, N. D. 2010a. Host  
593 association of *Campylobacter* genotypes transcends geographic variation. *Appl Environ*  
594 *Microbiol*, 76, 5269-77.

595 SHEPPARD, S. K., COLLES, F. M., MCCARTHY, N. D., STRACHAN, N. J., OGDEN, I. D., FORBES, K. J.,  
596 DALLAS, J. F. & MAIDEN, M. C. 2011. Niche segregation and genetic structure of  
597 *Campylobacter jejuni* populations from wild and agricultural host species. *Mol Ecol*, 20,  
598 3484-90.

599 SHEPPARD, S. K., DALLAS, J. F., STRACHAN, N. J., MACRAE, M., MCCARTHY, N. D., WILSON, D. J.,  
600 GORMLEY, F. J., FALUSH, D., OGDEN, I. D., MAIDEN, M. C. & FORBES, K. J. 2009.  
601 *Campylobacter* genotyping to determine the source of human infection. *Clin Infect Dis*, 48,  
602 1072-8.

603 SHEPPARD, S. K., DALLAS, J. F., WILSON, D. J., STRACHAN, N. J., MCCARTHY, N. D., JOLLEY, K. A.,  
604 COLLES, F. M., ROTARIU, O., OGDEN, I. D., FORBES, K. J. & MAIDEN, M. C. 2010b. Evolution of  
605 an agriculture-associated disease causing *Campylobacter coli* clade: evidence from national  
606 surveillance data in Scotland. *PLoS One*, 5, e15708.

607 SHEPPARD, S. K., DIDELOT, X., JOLLEY, K. A., DARLING, A. E., PASCOE, B., MERIC, G., KELLY, D. J.,  
608 CODY, A., COLLES, F. M., STRACHAN, N. J., OGDEN, I. D., FORBES, K., FRENCH, N. P., CARTER,  
609 P., MILLER, W. G., MCCARTHY, N. D., OWEN, R., LITRUP, E., EGHOLM, M., AFFOURTIT, J. P.,  
610 BENTLEY, S. D., PARKHILL, J., MAIDEN, M. C. & FALUSH, D. 2013a. Progressive genome-wide  
611 introgression in agricultural *Campylobacter coli*. *Mol Ecol*, 22, 1051-64.

612 SHEPPARD, S. K., DIDELOT, X., MERIC, G., TORRALBO, A., JOLLEY, K. A., KELLY, D. J., BENTLEY, S. D.,  
613 MAIDEN, M. C., PARKHILL, J. & FALUSH, D. 2013b. Genome-wide association study identifies  
614 vitamin B5 biosynthesis as a host specificity factor in *Campylobacter*. *Proc Natl Acad Sci U S*  
615 *A*, 110, 11923-7.

616 SHEPPARD, S. K., JOLLEY, K. A. & MAIDEN, M. C. 2012. A Gene-By-Gene Approach to Bacterial  
617 Population Genomics: Whole Genome MLST of *Campylobacter*. *Genes (Basel)*, 3, 261-77.

618 SHEPPARD, S. K., MCCARTHY, N. D., FALUSH, D. & MAIDEN, M. C. 2008. Convergence of  
619 *Campylobacter* species: implications for bacterial evolution. *Science*, 320, 237-9.

620 SMID, J. H., MUGHINI GRAS, L., DE BOER, A. G., FRENCH, N. P., HAVELAAR, A. H., WAGENAAR, J. A. &  
621 VAN PELT, W. 2013. Practicalities of using non-local or non-recent multilocus sequence  
622 typing data for source attribution in space and time of human campylobacteriosis. *PLoS One*,  
623 8, e55029.

624 SONG, W. S., NAM, M. S., NAMGUNG, B. & YOON, S. I. 2015. Structural analysis of PseH, the  
625 *Campylobacter jejuni* N-acetyltransferase involved in bacterial O-linked glycosylation.  
626 *Biochem Biophys Res Commun*, 458, 843-8.

627 TAMURA, K., STECHER, G., PETERSON, D., FILIPSKI, A. & KUMAR, S. 2013. MEGA6: Molecular  
628 Evolutionary Genetics Analysis version 6.0. *Mol Biol Evol*, 30, 2725-9.

- 629 VISWANATHAN, M., PEARL, D. L., TABOADA, E. N., PARMLEY, E. J., MUTSCHALL, S. & JARDINE, C. M.  
630 2016. Molecular and Statistical Analysis of *Campylobacter* spp. and Antimicrobial-Resistant  
631 *Campylobacter* Carriage in Wildlife and Livestock from Ontario Farms. *Zoonoses Public*  
632 *Health*.
- 633 WEBB, G. F. & BLASER, M. J. 2002. Dynamics of bacterial phenotype selection in a colonized host.  
634 *Proc Natl Acad Sci U S A*, 99, 3135-40.
- 635 WILSON, D. J., GABRIEL, E., LEATHERBARROW, A. J., CHEESBROUGH, J., GEE, S., BOLTON, E., FOX, A.,  
636 FEARNHEAD, P., HART, C. A. & DIGGLE, P. J. 2008. Tracing the source of campylobacteriosis.  
637 *PLoS Genet*, 4, e1000203.
- 638 WILSON, D. J., GABRIEL, E., LEATHERBARROW, A. J., CHEESBROUGH, J., GEE, S., BOLTON, E., FOX, A.,  
639 HART, C. A., DIGGLE, P. J. & FEARNHEAD, P. 2009. Rapid evolution and the importance of  
640 recombination to the gastroenteric pathogen *Campylobacter jejuni*. *Mol Biol Evol*, 26, 385-  
641 97.
- 642 WIMALARATHNA, H. M., RICHARDSON, J. F., LAWSON, A. J., ELSON, R., MELDRUM, R., LITTLE, C. L.,  
643 MAIDEN, M. C., MCCARTHY, N. D. & SHEPPARD, S. K. 2013. Widespread acquisition of  
644 antimicrobial resistance among *Campylobacter* isolates from UK retail poultry and evidence  
645 for clonal expansion of resistant lineages. *BMC Microbiol*, 13, 160.
- 646 YAHARA, K., DIDELOT, X., ANSARI, M. A., SHEPPARD, S. K. & FALUSH, D. 2014. Efficient inference of  
647 recombination hot regions in bacterial genomes. *Mol Biol Evol*, 31, 1593-605.
- 648 YAHARA, K., FURUTA, Y., OSHIMA, K., YOSHIDA, M., AZUMA, T., HATTORI, M., UCHIYAMA, I. &  
649 KOBAYASHI, I. 2013. Chromosome painting in silico in a bacterial species reveals fine  
650 population structure. *Molecular Biology and Evolution*, 30, 1454-64.

651

652

653 **Tables and Figures**

654 **Figure 1: Population structure of *Campylobacter* isolates used in this study.** Phylogenetic  
655 trees were constructed from a whole-genome alignment of **(A)** *C. jejuni* (n=229) and **(B)** *C.*  
656 *coli* (n=55) isolates based on 103,878 and 806,657 variable sites, respectively, using an  
657 approximation of the maximum likelihood algorithm (Tamura et al., 2013, Kumar et al.,  
658 2016). Leaves on the tree are coloured by source country, UK (green circles), Canada (red)  
659 and USA (blue). Common clonal complexes (CC) are annotated based on four or more shared  
660 alleles in seven MLST house-keeping genes (Dingle et al., 2005).

661

662 **Figure 2: Co-ancestry matrix with population structure and genetic flux.** **(A)** The colour  
663 of each cell of the matrix indicates the number of chunks imported from a donor genome  
664 (column) to a recipient genome (row). Colour ranges from little import from the donor strain  
665 (yellow) to a large amount of imported DNA from the donor strain (blue). White indicates  
666 missing data. The trees above and to the left show clustering of the paired isolates with leaves  
667 coloured by source country (UK in green, Canada in red). **(B)** Box plot summarising the co-  
668 ancestry matrix data. The average frequency of inferred recombination between donor to  
669 recipient grouped by import from isolates from the same country compared to isolates from  
670 different countries. There is significantly more import from donor strains of the same country  
671 compared to strains from different countries (*p-value, test*).

672

673 **Figure 3: Pairwise comparison of nucleotide diversity in the core genome.** **Above:**  
674 Estimated values of the per-nucleotide statistic reflecting relative intensity of recombination  
675 at each site plotted along the NCTC11168 reference genome. **Left:** Core genome phylogeny  
676 of selected paired isolates (matched by CC and source host), with clonal complex indicated.

Commented [XD16]: To insert

677 **Centre:** Matrix of gene-by-gene pairwise comparison along the NCTC11168 reference  
678 genome of our selected pairs. Each row represents a pairwise comparison of selected paired  
679 of isolates. Each column is a gene from the NCTC11168 reference genome. Panels of the  
680 matrix are coloured based on nucleotide divergence for that gene in each pair: from no  
681 nucleotide diversity (0%, white), through some nucleotide diversity (~1%, red) to high levels  
682 of nucleotide diversity (up to 2%, blue). The per-nucleotide scan of relative intensity of  
683 recombination is aligned with our gene-by-gene pairwise comparison of nucleotide diversity  
684 and the location of **seven** putative epidemiological markers for geographical segregation are  
685 indicated.

**Commented [XD17]:** I can only count six

686  
687 **Figure 4: Assignment of human clinical cases of campylobacteriosis to origin country,**  
688 **including patients with history of recent foreign travel.** (A) Assignment of human clinical  
689 cases of campylobacteriosis to origin country using epidemiological markers of biogeography  
690 and the Bayesian clustering algorithm Structure. Each isolate is represented by a vertical bar,  
691 showing the estimated probability that it comes from each of the putative source countries,  
692 including the UK (green), **USA (blue)** and Canada (red). Isolates are ordered by attributed  
693 source. (B) Boxplots of predicted attribution probabilities for the three locations. (C) Isolates  
694 from Oxford clinical dataset with declared history of recent foreign travel. The model  
695 correctly assigned 34 of 46 (73.9%) isolates to a non-UK origin. (D) Attribution of Oxford  
696 clinical isolates between UK, USA and Canada source populations. Isolates with declared  
697 recent foreign travel are shown in blue.

**Commented [ELS18]:** The key on the figure needs to be changed from yellow to blue

698  
699 **Table 1:** Isolate pairs matched by clonal complex and host.

700

701 **Table 2: Shared ancestry analysis and estimation of pairwise recombination rates.** The  
702 time to the most recent common ancestor (TMRCAs) for each selected pair was estimated  
703 with 95% confidence intervals (TMRCAs-CI). The ratio of rates at which recombination and  
704 mutation introduce polymorphism ( $r/m$ ) was also calculated with 95% confidence intervals  
705 ( $r/m$ -CI). In addition, the number of definitely recombined genes (probability > 95%) is also  
706 shown. The two *C. coli* pairs are coloured in red.

707

#### 708 **Supplementary material**

709 **Figure S1:** Phylogeny of 7 highly recombining epidemiological markers used to attribute  
710 biogeography using structure.

711 **Table S1:** List of isolates used, including details of genome accession numbers.

712 **Table S2:** List of biogeographical epidemiological markers, including lists of (A) highly  
713 recombining genes as determined by per-nucleotide estimation of recombination intensity  
714 (recombination hot spots); (B) highly recombining genes as determined by pairwise analysis  
715 of nucleotide diversity (more than 2% diversity); and genes used to model biogeographical  
716 segregation in structure (orange). Genes with a role in fluoroquinolone resistance are  
717 highlighted in yellow.

718