

1 **Quantifying bacterial evolution in the wild: a birthday problem for *Campylobacter***  
2 **lineages**

3  
4 Jessica K. Calland<sup>1</sup>, Ben Pascoe<sup>1</sup>, Sion C. Bayliss<sup>1</sup>, Evangelos Mourkas<sup>1</sup>, Elvire Berthenet<sup>2</sup>,  
5 Harry A. Thorpe<sup>1</sup>, Matthew D. Hitchings<sup>3</sup>, Edward J. Feil<sup>1</sup>, Martin J. Blaser<sup>4</sup>, Daniel Falush<sup>5</sup>#  
6 & Samuel K. Sheppard<sup>1</sup>#

7  
8 <sup>1</sup>The Milner Centre for Evolution, University of Bath, Claverton down, Bath, UK; <sup>2</sup>French  
9 National Reference Center for Campylobacters and Helicobacters, University of Bordeaux,  
10 33076, Bordeaux, France; <sup>3</sup>Institute of Life Sciences, Swansea University Medical School,  
11 Swansea University, Singleton Park, Swansea, UK; <sup>4</sup>Center for Advanced Biotechnology and  
12 Medicine, Rutgers University, New Brunswick, New Jersey, USA; <sup>6</sup>Centre for Microbes,  
13 Development and Health, Institute Pasteur Shanghai, China.

14  
15 #Authors to which correspondence should be addressed: [s.k.sheppard@bath.ac.uk](mailto:s.k.sheppard@bath.ac.uk) &  
16 [daniel.falush@ips.ac.cn](mailto:daniel.falush@ips.ac.cn)

17  
18 **Short title:** *Campylobacter* molecular clock

19  
20 **Keywords:** Evolution / Pathogenesis / Molecular clock / Population structure / Population  
21 genomics.

22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33 **Abstract**

34 Measuring molecular evolution in bacteria typically requires estimation of the rate at which  
35 mutations accumulate in strains sampled at different times that share a common ancestor. This  
36 approach has been useful for dating ecological and evolutionary events that coincide with the  
37 emergence of important lineages, such as outbreak strains and obligate human pathogens.  
38 However, in multi-host (niche) transmission scenarios, where the pathogen is essentially an  
39 opportunistic environmental organism, sampling is often sporadic and rarely reflects the overall  
40 population, particularly when concentrated on clinical isolates. This means that approaches that  
41 assume recent common ancestry are not applicable. Here we present a new approach to estimate  
42 the molecular clock rate in *Campylobacter* that draws on the popular probability conundrum  
43 known as the ‘birthday problem’. Using large genomic datasets and comparative genomic  
44 approaches, we identify isolate pairs where common ancestry is inferred within the sample  
45 time-frame – analogous to a shared birthday. Identifying synonymous and non-synonymous  
46 substitutions, both within and outside of recombinant regions of the genome, we quantify clock-  
47 like diversification to estimate mutation rates for the common pathogenic species  
48 *Campylobacter coli* ( $2.4 \times 10^{-6}$  s/s/y) and *Campylobacter jejuni* ( $3.4 \times 10^{-6}$  s/s/y). Finally, using  
49 estimated mutation rates we assess the rate of turnover of lineages in our sample set over short  
50 evolutionary timescales. This provides a generalizable approach to calibrating mutation rates in  
51 populations of environmental bacteria and shows that multiple lineages are maintained,  
52 implying that large-scale clonal sweeps may take hundreds of years or more in these species.

53

54

55

56

57

58

59

60

61

62

63

64

## 65 **Introduction**

66 Theoretical models of a relatively constant rate of molecular change over time (Kimura 1968),  
67 the molecular clock, have become fundamental to explaining the evolution in bacteria (Kuo and  
68 Ochman 2009; Didelot et al. 2016). Spurred by the increasing availability of population-scale  
69 genome datasets, it is now common for comparative genomic studies to describe not only the  
70 relatedness of isolates but also how long ago they diverged (Kidgell et al. 2002; Mutreja et al.  
71 2011; Mcadam et al. 2012; Cui et al. 2013; Mourkas et al. 2020). This can provide valuable  
72 information when combined with host, habitat or ecosystem data. For example, it is possible to  
73 investigate how events such as host transitions or global dissemination have influenced the  
74 emergence and spread of lineages that may display important phenotypes, including  
75 pathogenicity.

76

77 There are significant challenges when applying molecular clocks to date lineage diversification  
78 in natural bacterial populations. In particular, it is necessary to determine the rate at which the  
79 clock ‘ticks’ and the uniform accumulation of nucleotide substitutions over time. However, this  
80 is not simply a reflection of the background point mutation rate (associated with replication  
81 error) and the generation time of the bacterium (Weller and Wu 2015; Gibson et al. 2018), but  
82 is also influenced by horizontal gene transfer (HGT) that can introduce several mutations in a  
83 single event (Vos and Didelot 2009). Furthermore, the rate at which mutations accumulate in  
84 the population is influenced by the population size (Bromham 2009) and selection (positive and  
85 stabilizing) on different fitness effects (Eyre-Walker and Keightley 2007).

86

87 While debate continues about mutations that are effectively neutral, and hence provide accurate  
88 clock estimates (Gibson and Eyre-Walker 2019), there is clear utility for even approximations  
89 of the rate of genome change over time (Drummond et al. 2003; Biek et al. 2015). This has  
90 allowed the development of time-calibrated phylogenies explaining molecular evolution in  
91 numerous well-known pathogen species (Kidgell et al. 2002; Mutreja et al. 2011; Mcadam et  
92 al. 2012; Cui et al. 2013). However, even with large genome datasets and increasingly  
93 sophisticated models (Drummond and Rambaut 2007; Suchard et al. 2018), the accuracy of  
94 molecular evolution estimates is dependent upon the data from which they are derived, and two  
95 important considerations remain. First, the data should represent a longitudinal sample set

96 (Drummond et al. 2003; Arnold and Hanage 2017). Second, the data should be representative  
97 of the population as a whole.

98

99 It is conceptually simple to understand how a long time-frame between collection of the earliest  
100 and latest sample would increase the number of mutations recorded, and how sampling at  
101 consistent intervals could help to determine if accumulation was linear over time. Comparisons  
102 between modern samples and DNA from the stomach of a 5,300 year old frozen iceman ‘Otzi’  
103 have been used to investigate the emergence of modern *Helicobacter pylori* lineages (Maixner  
104 et al. 2016). However, ancient pathogen samples are rarely available. More frequently,  
105 molecular clock rates are estimated using collections of contemporary isolates that often share  
106 a common ancestor older than the sample frame. Convincing estimations have been possible  
107 for medically important bacteria, through comparison of large numbers of closely related  
108 isolates (Didelot et al. 2012; Walker et al. 2013; Mathers et al. 2015; Menardo et al. 2019) but  
109 for many pathogens sampling of outbreaks may not provide an adequate representation of the  
110 bacterial population.

111

112 Most disease-causing bacteria are not obligate human pathogens. In this case, large reservoirs  
113 of isolates from which infection can arise may be infrequently sampled, despite their potential  
114 importance as emergent pathogenic strains. For example, *Campylobacter jejuni* and *C. coli* are  
115 among the most common causes of bacterial gastroenteritis worldwide but exist principally as  
116 commensal organisms in the gut of mammals and birds (Waldenström et al. 2002; Sheppard et  
117 al. 2011; Bronowski et al. 2014; Cody et al. 2015; Sheppard and Maiden 2015). Human  
118 infection results primarily via food contaminated with strains from wild and agricultural  
119 animals, especially chickens (Wilson et al. 2008; Sheppard et al. 2009; Strachan et al. 2009;  
120 Dearlove et al. 2016; Rosner et al. 2017; Thépault et al. 2017). In multi-host (niche)  
121 transmission scenarios such as this, where the pathogen is essentially an environmental  
122 organism, sampling is often sporadic and rarely reflects the overall population, particularly  
123 when concentrated on clinical isolates (Marin and Hedges 2018).

124

125 Overcoming the problem of sporadic or unrepresentative sampling for molecular clock  
126 estimation requires that sufficient numbers of isolates are collected to ensure that there are pairs  
127 that share a recent common ancestor (within the sampling period). However, with the enormous



128 effective population size of environmental bacteria populations, questions remain about how  
129 many isolates need to be sampled to achieve this. This is analogous to the well-known  
130 probability theory conundrum known as the birthday problem (Mathis 1991). This puzzle asks  
131 how many randomly chosen people need to be sampled so that a pair of them will share the  
132 same birthday. To be sure, requires a sample size of 366 (the number of possible birthdays),  
133 assuming that all birthdays are equally common, but a 99.9% probability is achieved with just  
134 70 people and 50% with 23 people. This may seem counter intuitive but can be explained by  
135 considering that rather than comparing the birthday of a single individual to everyone else's, in  
136 fact comparisons are made between every pair of individuals,  $23 \times 22/2 = 253$ . The result is  
137 greater than half the number days in the year, hence the 50% probability. Clearly, there are  
138 challenges in relating this conceptual model to bacteria. First, it is not known how many  
139 possible lineages (here equivalent to birthdays) there are in natural bacterial populations.  
140 Second, how to define lineages or isolate pairs with recent common ancestry. Third, just as with  
141 birthdays, some lineages are far more common than others. For example, of >72,000 *C. jejuni*  
142 and *C. coli* isolates archived in the pubMLST database (Jolley et al. 2018), >50% belong to just  
143 5 clonal complexes (out of 45).

144

145 Together, factors relating to isolate sampling and genome analysis conspire such that it may be  
146 difficult to distinguish nucleotide substitutions that reflect the passage of time (Didelot and  
147 Falush 2007; Biek et al. 2015). Here, we take a multi-layered approach to estimate the rate of  
148 molecular evolution of *C. jejuni* and *C. coli* using a large genome collection (2,425 genomes)  
149 representing isolates sampled over a 46-year period. We begin by identifying isolate pairs  
150 where: (i) common ancestry is inferred within the sample time-frame, and (ii) the most recently  
151 sampled isolate has accumulated SNPs over time. We then quantify synonymous and non-  
152 synonymous polymorphisms to take (some) account of selection, both within and outside of  
153 recombinant regions of the genome, and use synonymous polymorphisms to quantify clock-  
154 like diversification in *Campylobacter* (Kimura 1987; Gojobori et al. 1990). Finally, using  
155 estimated mutation rates we assess the rate of turnover of lineages in our sample sets over short  
156 evolutionary timescales. This provides a generalizable approach to calibrating mutation rates in  
157 populations of environmental bacteria and clues about lineage diversification in two important  
158 pathogenic bacteria.

159

## 160 **Results**

### 161 ***There is a weak temporal signal in C. coli and C. jejuni phylogenies***

162 Core genome phylogenies revealed little evidence of clustering by collection date (**Figure 1**).  
163 Isolates belonging to common sequence types (STs) and clonal complexes were sampled over  
164 the 46-year period. These included poultry associated ST-353, ST-354 and ST-257 complexes,  
165 cattle associated ST-61 and ST-42 complexes, and host generalist ST-21, ST-45, ST-828 (*C.*  
166 *coli*) complexes (Sheppard et al. 2014) (**Figure 1 and Supplementary Table 1**). Linear  
167 regression of root-to-tip distances and sampling dates of *C. coli* and *C. jejuni* phylogenies  
168 (**Supplementary Figure 1**), using TempEst software, provided very weak evidence of a  
169 temporal signal when the best-fitting root was estimated. The  $R^2$  values were low for both *C.*  
170 *coli* ( $R^2 = 0.176$ ) and *C. jejuni* ( $R^2 = 9.5 \times 10^{-2}$ ) phylogenies (**Supplementary Table 2**).  
171 Consistent with some other studies (Rieux and Balloux 2016), this poor branch-length to  
172 isolation date correlation suggests that estimation of the molecular clock rate from the entire  
173 dataset may be difficult. However, the accumulation of polymorphisms exhibited a positive  
174 correlation with sampling date in all datasets (**Supplementary Figure 1**) implying the  
175 maintenance of multiple STs and clonal complexes through time.

176

### 177 ***Sampling matched isolate pairs allows estimation of mutation rate***

178 Estimation of molecular clock rates requires comparison of isolates from related, or preferably  
179 the same lineages, that have accumulated mutations over time. To achieve this there is a  
180 necessary balance between maximizing the time between sampling and accumulated SNPs  
181 whilst ensuring comparisons are made between related strains. Therefore, we plotted SNP  
182 difference against time difference to determine criteria for choosing comparable isolate pairs  
183 (**Figure 2**). The sample time difference was chosen to maximize the time between sampling  
184 and the number of comparable pairs belonging to the same lineage. Pair selection criteria were  
185 standardised for both species so that isolate pairs were excluded where the sampling time  
186 difference was <8 years or there were >5000 SNPs between them (**Figure 2, Supplementary**  
187 **Table 3, Table 1**). Based upon these criteria, 18 *C. coli* and 20 *C. jejuni* isolate pairs were  
188 considered for mutation rate calibration. These belonged to the ST-21, ST-22, ST-45, ST-1332,  
189 ST-828 clonal complexes and isolate pairs had a difference in sampling date of 8 to 11 years  
190 (*C. coli*) and 8 to 36 years (*C. jejuni*) (**Tables 1 and 2, Supplementary Table 3**).

191

192 Estimation of a molecular clock rate requires that SNPs accumulate over time, defined here as  
193 mutations per site per year (s/s/y). It is also possible that branch shortening can occur where  
194 there are fewer mutations in the more recent isolate of a pair resulting in a negative rate of  
195 molecular evolution (Duchêne et al. 2016). In this study, 13 out of 18 *C. coli* and 11 out of 20 *C.*  
196 *jejuni* isolate pairs exhibited branch lengthening, i.e. more total mutations (within and outside  
197 recombined regions) were found in the more recent isolate (**Tables 3 and 4, Supplementary**  
198 **Table 4**). Only pairs having undergone measurable evolution (branch lengthening) were  
199 included in further analysis of the accumulation of mutations over time. For these isolate pairs,  
200 the total mutation rate was calculated as well as the mutation rate within and outside of  
201 recombined regions (**Table 5, Supplementary Table 5**). The mean mutation rate for non-  
202 recombined regions was  $6.36 \times 10^{-6}$  and  $8.45 \times 10^{-6}$  s/s/y for *C. coli* and *C. jejuni* respectively,  
203 or 11.46 and 13.53 mutations per genome per year (s/g/y) (**Table 5**).

204

#### 205 ***Recombination drives molecular evolution in Campylobacter***

206 Mutations in coding sequence based on gene definitions in reference isolate genomes (CVM  
207 N29710 (*C. coli*) and NCTC 11168 (*C. jejuni*)) introduced an average of 1569 and 242 SNPs  
208 in *C. coli* and *C. jejuni* paired genome datasets respectively. Of these, an average of only 222  
209 (*C. coli*) and 106 (*C. jejuni*) were inferred to be the result of point mutation, with the remainder  
210 resulting from recombination (**Tables 3 and 4**). Recombination is therefore the major source  
211 of sequence variation in both species (**Figure 3, Tables 3 and 4**), introducing nearly six times  
212 as many polymorphisms in *C. coli* than in *C. jejuni* – consistent with previous estimates based  
213 upon MLST (Wilson et al. 2009). To assess the effect of mutations on amino acid sequences  
214 we quantified non-synonymous (N) and synonymous (S) mutations and determined the ratio  
215 per site ( $dN/dS$ ) for all isolate pairs in recombined and non-recombined sequence (**Tables 3 and**  
216 **4**). Point mutation accounted for an unequal amount of N and S polymorphism both within and  
217 between species (*C. coli*, N = 99, S = 123; *C. jejuni*, N = 63, S = 43). While recombination  
218 introduced many more polymorphisms than point mutation, in both species these were biased  
219 towards synonymous changes. Specifically, around six times as many S than N mutations were  
220 introduced by recombination in *C. coli* and approximately twice more in *C. jejuni* (*C. coli*, N =  
221 546, S = 801; *C. jejuni*, N = 59, S = 77). Overall, average  $dN/dS$  ratios were consistent between  
222 species within recombined (*C. coli* 0.492, *C. jejuni* 0.490) and non-recombined (*C. coli* 0.594,  
223 *C. jejuni* 0.509) portions of the genome. However, because of the relative importance of

224 recombination ( $r/m = 37.240$  (*C. coli*),  $r/m = 5.098$  (*C. jejuni*)), on average N mutations were  
225 similar for *C. jejuni* from recombination and point mutation (59 and 63 respectively). However,  
226 recombination introduced 5.5 times more N mutations than point mutation in *C. coli* (**Tables 3**  
227 **and 4**). Variation in  $dN/dS$  was observed between isolate pairs but was mostly indicative  
228 purifying selection ( $dN/dS < 1$ ). Evidence of positive selection ( $dN/dS > 1$ ) was only observed  
229 within recombined sequence in 6 isolate pairs (**Tables 3 and 4**).

230

231 Additional analysis of the distribution of recombination events revealed that an average of 13%  
232 (*C. coli*) and 2% (*C. jejuni*) of the genome has undergone recombination in at least one isolate  
233 pair since divergence from the common ancestor of each sub-tree (**Tables 6 and 7**).  
234 Recombination was distributed across the genome in both species but was elevated in certain  
235 regions of *C. coli* introducing more polymorphism at potential recombination hotspots (Yahara  
236 et al. 2014). However, recombination remained the main source of variation in both species  
237 (**Figure 3**).

238

### 239 *Molecular clock estimates for C. coli and C. jejuni*

240 Molecular clock estimates require that mutations accumulate at a consistent rate over time. We  
241 maximized the chance of identifying this signal in several ways. First, genomic variation within  
242 recombined regions was discounted as multiple SNPs can be introduced in a single evolutionary  
243 event – distorting clock estimates (Didelot and Falush 2007; Wilson et al. 2009; Croucher et al.  
244 2011). Second, non-synonymous mutations were discounted as selection may be more likely to  
245 influence the frequency of variation at these sites. Third, only pairs in which the most recently  
246 sampled isolate contained more SNPs (branch lengthening) were used as they displayed  
247 measurable evolution over time. Based on these criteria, a similar average molecular clock rate  
248 was obtained for *C. coli*,  $2.4 \times 10^{-6}$  s/s/y (4.27 s/g/y), and *C. jejuni*,  $3.4 \times 10^{-6}$  s/s/y (5.42 s/g/y)  
249 (**Table 5**).

250

### 251 *Coalescence and maintenance of lineages over time*

252 Molecular clock estimates can vary within a population. Therefore the applicability of  
253 generalized clocks depend upon how much of the population has been sampled. To quantify  
254 this we estimated the average mutation rate ( $\mu$ ) (*C. coli* = 77.292 s/g/y, *C. jejuni* = 14.101 s/g/y),  
255 including all polymorphisms within and outside recombined sequence. These mutation rates

256 were used to determine the number of coalescences in the population at a given time point (here  
257 referred to as ‘*effective lineages*’) within the dataset. The maximum timeframe for comparison  
258 was 37 years for *C. coli* and 46 years for *C. jejuni* (short in evolutionary terms). This provided  
259 information about the number of ancestral strains and the rate of turnover of lineages within the  
260 dataset. The total number of potential pairs without accounting for genetic similarity ( $Y$ ), was  
261 equal to the square of the total number of isolates ( $n^2$ ) divided by two (to avoid double counting  
262 of isolate pairs), 180,600 and 1,663,488 for *C. coli* and *C. jejuni* respectively.

263  
264 Having determined the mutation rate, we were able to predict the expected number of mutations  
265 over a given period of time. For example, 14 in 1 year for *C. jejuni*. We then subsampled all  
266 isolate pairs ( $Y$ ) to determine how many isolate pairs had  $\leq 14$  SNPs between them – 76 isolate  
267 pairs. This is the possible number of isolate pairs that have arisen in 1 year. This process was  
268 repeated for each time cut-off, up to a maximum of 37 and 46 years for *C. coli* and *C. jejuni*  
269 respectively (**Tables 8 and 9**), to give the number of possible pairs for every time cut-off ( $X$ )  
270 (**Figure 2B and 2D**). Dividing  $Y/X$  resulted in the number of coalescences (*effective lineages*)  
271 at a given time interval in the past ( $Z$ ) (**Tables 8 and 9**). For example, if the total mutation rate  
272 was 14 s/g/y and we were interested in the number of birthdays within 5 years of our dataset,  
273 we would multiply the mutation rate by 5 to result in 70 SNPs of evolution over 5 years. The  
274 number of *potential pairs* ( $Y = 1,663,488$ ) / *possible pairs* ( $X = 174$ ) = ~9,560 coalescences  
275 (ancestors) within this time period (**Supplementary Figure 2, Tables 8 and 9**).

276  
277 The number of effective lineages at a given time-point can also be interpreted as the number of  
278 lineages that gave rise to those that are seen today. This provides valuable information about  
279 how the population is maintained over time and the extent to which it has diversified. For  
280 example, 1,263 *C. coli* lineages 37 years ago gave rise to an estimated 22,575 one year ago and  
281 4,726 *C. jejuni* lineages 46 years ago gave rise to 21,888 lineages one year ago. This equates to  
282 an average increase in the number of effective lineages of 576 and 373 per year for *C. coli* and  
283 *C. jejuni* respectively. For *C. jejuni* it is clear that a considerable proportion (22%) of all  
284 lineages have been maintained throughout the 46 year sampling period and probably much  
285 longer (**Figure 4**). In contrast, only 6% of all effective lineages were present in the *C. coli*  
286 population 37 years ago. Perhaps the most striking finding is that the *C. coli* population has

287 rapidly diversified in recent years. For example, there has been an 800% increase in the number  
288 of effective lineages in the last 10 years, over 3 times the rate of increase observed in *C. jejuni*.  
289 **(Figure 4, Supplementary Figure 2).**

290

## 291 **Discussion**

292 The increasing availability of large genome datasets has great potential for improving molecular  
293 clock estimates in bacteria. However, significant challenges remain. While it is clear that the  
294 frequency of substitutions can vary between different species and strains (von Mering et al.  
295 2007; Mcadam et al. 2012; Cui et al. 2013; Li et al. 2015; Duchêne et al. 2016; Gibson et al.  
296 2018; Menardo et al. 2019), the extent to which nucleotide variation represents an intrinsic  
297 molecular clock is often less apparent. Biological factors such as generation time, population  
298 size and recombination rate, and ecological factors including cellular responses to habitat  
299 variation or stress and the strength of natural selection, influence the rate at which substitutions  
300 accumulate in populations (Denamur and Matic 2006). Therefore, obtaining a robust molecular  
301 clock estimate from natural bacterial populations requires an appropriate sample frame and  
302 careful consideration of the nature of observed sequence variation.

303

304 In cases where there is a clear temporal signal among isolates, it may be possible to obtain a  
305 robust molecular clock estimate by applying models to large genome datasets (Menardo et al.  
306 2019). However, analysing all *C. jejuni* and *C. coli* genomes together gave a weak temporal  
307 signal. This is likely related to the population structure and biology of these organisms that is  
308 in stark contrast to many obligate human pathogens (Menardo et al. 2019). Consistent with  
309 many other zoonotic or environmental bacteria, *Campylobacter* is a diverse genus with multiple  
310 lineages (STs and clonal complexes) inhabiting multiple hosts/niches. This required a more  
311 targeted approach to microevolutionary analysis consistent with that used to investigate  
312 transmission in similarly variable organisms (Didelot et al. 2012).

313

314 Sub-sampling within the isolate collection, sampled over 46 years, identified closely related  
315 pairs of isolates with divergent sampling dates. Clearly, calibration of the molecular clock  
316 requires that mutations accumulate over time. This was not the case in all isolate pairs. In some  
317 cases, the most recently sampled isolate had accumulated fewer substitutions than the  
318 comparator strain leading to a negative mutation rate as observed in some other bacterial species



319 (Duchêne et al. 2016; Menardo et al. 2019). This could indicate time-dependency of molecular  
320 evolution (Ho et al. 2007; Ho et al. 2011), where deleterious mutations in the older isolate have  
321 been purged leading to differences in long and short term molecular clock estimates (Rocha et  
322 al. 2006; Duchêne et al. 2014). However, in organisms with complex ecology such as  
323 *Campylobacter*, it is also possible that closely related isolates occupy different sub-niches and  
324 experience different selection pressures even when sampled from the same host.

325

326 Returning to the birthday problem analogy, considering the number of isolate pairs (equivalent  
327 to people with the same birthday) obtained from the original genome dataset can provide clues  
328 about the extent of lineage diversity in the natural population. Using total mutation rates, we  
329 were able to assess the nature of coalescence across the sample time frame for each species.  
330 The coalescence we refer to here is equivalent to the number of ancestral strains at a particular  
331 time point (effective lineages) in the natural environment from which contemporary strains  
332 emerged. Effective population size ( $N_e$ ) is commonly used to reflect the number of individuals  
333 in a population that contribute to subsequent generations (Kirchberger et al. 2020). This has  
334 been used to investigate bacteria but contrasting approaches can provide different estimates  
335 depending on the method used (Cui et al. 2015; von Mering et al 2007). The idea of effective  
336 lineages, described in this study, is related to  $N_e$  but is more specific for clonal organisms.  
337 Rather than typical  $N_e$  estimates for sexual populations, where the mating of two individuals is  
338 largely independent of what happened in previous generations, the number of effective lineages  
339 in a bacterial population reflects the number of distinct lineages that will survive and therefore  
340 contribute to future generations. This provides information on the genetic inertia of the  
341 population.

342

343 These analyses highlighted the importance of appropriate sampling when calibrating mutation  
344 rates and can help in determining the extent to which samples represent the population as a  
345 whole. Specifically, by considering the number of coalescences in a random population, we can  
346 look back through the sample time frame to estimate the number of effective lineages across a  
347 randomly sampled dataset. For example, suppose we would like to know if our contemporary  
348 isolates have a common ancestor in 1980. We know that a proportion of these ancestors gave  
349 rise to the diversity we see today but many lineages would go extinct and therefore not  
350 contribute (Louca et al. 2018). Based on an average mutation rate of 14 s/g/y for *C. jejuni*,

351 there would be 560 SNPs over 40 years total evolution between a strain pair. So, one can then  
352 ask how many pairs are close enough genetically for that to be the case. This gives an estimate  
353 of the effective number of ancestors in 1980 that gave rise to the contemporary dataset -  
354 equivalent to the number of birthdays.

355

356 For *Campylobacter*, it is clear that multiple lineages have persisted over a long period of time.  
357 This indicates that although the population is large, the strains are not turning over at a  
358 particularly fast rate and are maintained over time. The absence of lineage replacement is  
359 inconsistent with some models of bacterial evolution that predict periodic population  
360 bottlenecks (Koeppel et al. 2007) but this can be explained in several ways. First, it is possible  
361 that the 37/46 year sampling period in this study is not sufficient time to out-compete a rival  
362 strain. Second, bacteria occupy different niches that are sustained so strains are not in direct  
363 competition. Third, the fitness differences among strains are not great enough for one lineage  
364 to out-compete another.

365

366 As well as the maintenance of multiple lineages, there is also evidence for variation in the  
367 number of effective lineages that contributed to successive generations between the two major  
368 pathogenic *Campylobacter* species. While this was consistently higher for *C. jejuni* throughout  
369 much of the sample frame there was a rapid increase in the number of *C. coli* lineages that  
370 began around 8 years ago (**Figure 4**). The reason for this is unclear. The average mutation rate  
371 estimates were similar for *C. jejuni* and *C. coli*,  $3.4 \times 10^{-6}$  and  $2.4 \times 10^{-6}$  s/s/y respectively,  
372 equating to approximately 5.4 (*C. jejuni*) and 4.3 (*C. coli*) mutations per genome per year. This  
373 is somewhat lower than previous estimates for *C. jejuni* calculated from 7-locus MLST ( $2.79 \times$   
374  $10^{-5}$  s/s/y) (Wilson et al. 2009) but is within the range of molecular clock estimates calculated  
375 from genomic variation for *Enterococcus faecium* ( $9.35 \times 10^{-6}$  s/s/y) *Y. pestis* ( $1.57 \times 10^{-8}$  s/s/y)  
376 (Duchêne et al. 2016).

377

378 While the average mutation rate was consistent for *C. coli* and *C. jejuni*, the relative number of  
379 polymorphisms introduced by homologous recombination and mutation ( $r/m$ ) differed  
380 markedly, with on average 37-fold (*C. coli*), compared to 5-fold (*C. jejuni*), greater impact on  
381 sequence variation. HGT is known to be an important driver of genome evolution in  
382 *Campylobacter* (Wilson et al. 2009; Sheppard et al. 2010) but these estimates are considerably



383 higher than previous ones using 7-locus MLST (Vos and Didelot 2009). Recombination  
384 introduced nearly twice as many synonymous than non-synonymous mutations, but even taking  
385 this into account, recombined sequence accounted for around 79% of all non-synonymous  
386 variation. This highlights the importance of HGT in rapidly evolving *Campylobacter* genomes  
387 and provides evidence that recombination may have been an important factor in the recent  
388 diversification of *C. coli* (Sheppard et al. 2008; Sheppard et al. 2011; Sheppard et al. 2013),  
389 potentially associated with an adaptive radiation (Rainey and Travisano 1998; Flohr et al. 2013)  
390 linked to the colonization of agricultural niches (Thakur et al. 2006). However, this should be  
391 balanced against the evidence of purifying selection within recombined sequence ( $dN/dS =$   
392 0.492 for *C. coli* and 0.49 for *C. jejuni*) and the removal of non-synonymous mutations through  
393 negative selection (Rocha et al. 2006).

394

395 Finally, throughout this study we have emphasized the importance of sampling so that measures  
396 of molecular evolution are obtained by comparing recent samples with a true ancestor. The  
397 uneven distribution of lineages within the population and the possibility that they differ in key  
398 evolutionary measures ( $r/m$  and  $dN/dS$ ), means that our molecular clock estimate may not be  
399 applicable to all *Campylobacter* lineages (Didelot et al. 2012; Croucher et al. 2013; Didelot et  
400 al. 2013; Everitt et al. 2014). Perhaps this is best illustrated by considering two host-specialist  
401 *C. jejuni* lineages, one associated with chickens and the other with cattle (Sheppard et al. 2011;  
402 Mourkas et al. 2020). There are 19 billion chickens on earth compared to 1.3 billion cattle (Bar-  
403 On et al. 2018) and *C. jejuni* colonizes up to 80% of chickens (Dhillon et al. 2006) with much  
404 lower rates in cattle. As the efficiency by which natural selection acts on sequence variation is  
405 related to effective population size (Gojobori et al. 1990), the rate of fixation and removal of  
406 mutations will be much faster in *C. jejuni* in chickens. Furthermore, chickens have a higher  
407 body temperature than cattle therefore the *C. jejuni* will grow faster, have a shorter generation  
408 time, and accumulate mutations at a higher rate (Weller and Wu 2015). From this simple  
409 example, which ignores many important factors (eg. subniche structure, host transition bottle  
410 necking, resident microbiome) it is clear molecular evolution can be influenced by population-  
411 scale forces down to the physiology of the individual cell. The approach employed in this study  
412 goes some way towards mitigating effects that confound generalized molecular clock estimates.  
413 Focussing on well-defined closely related isolate pairs inevitably reduces the number of  
414 comparisons from which the mean molecular clock rate is estimated. However, consideration

415 of the distribution of effective lineages within the population is essential for identifying robust  
416 molecular clock estimates in environmental bacteria with complex multi-host ecology and  
417 massive effective population sizes.

418

## 419 **Materials and Methods**

### 420 ***Isolate sampling, genome sequencing and assembly***

421 The accuracy of molecular clock estimates are improved by sampling strains over long time  
422 periods. To achieve this, an isolate collection was assembled comprising 53 isolates sampled  
423 between 1978 and 1985 (12 *C. coli*, 41 *C. jejuni*) and derived from multiple sources (human,  
424 duck, cattle, dog, turkey, wild bird and pig (**Supplementary Table 1**)). Samples were streaked  
425 onto mCCDA (PO0119A Oxoid Ltd, Basingstoke, UK) with CCDA Selective Supplement  
426 (SR0155E Oxoid Ltd, Basingstoke, UK) and incubated at 37°C for 48h in a microaerobic  
427 atmosphere (85% N<sub>2</sub>, 10% CO<sub>2</sub>, and 5% O<sub>2</sub>) using CampyGen Compact sachets (Thermo Fisher  
428 Scientific Oxoid Ltd, Basingstoke UK). Single colonies from each plate was then sub-cultured  
429 onto Mueller Hinton (MH) (CM0337 Oxoid Ltd, Basingstoke, UK) agar and grown for an  
430 additional 48h at 37°C and stored in 20% glycerol stocks at -80°C.

431

432 DNA was extracted using the QIAamp DNA Mini Kit (QIAGEN, Crawley, UK), according to  
433 manufacturer's instructions. DNA was quantified using a Nanodrop spectrophotometer before  
434 sequencing on an Illumina MiSeq sequencer using the Nextera XT library preparation kits with  
435 standard protocols. Paired end libraries were sequenced using 2 × 300 bp 3rd generation reagent  
436 kits (Illumina). Short read data was assembled using the *de novo* assembly algorithm, SPAdes  
437 (version 3.10.0 35) (Bankevich et al. 2012) generating an average of 49 contigs (range: 2 -115)  
438 for a total average assembled genome size of 1.69 Mbp (range: 1.62-1.80). The average N50  
439 was 189,430 bp (range: 81,283-974,529). These isolate genomes were augmented with 1,783  
440 *C. jejuni* and 589 *C. coli* genomes archived in BIGSdb (Jolley and Maiden 2010) representing  
441 isolates sampled from multiple sources (human, cattle, chicken, cat, dog, duck, environmental  
442 waters, farm environments, geese, lamb, rabbit, sand, seal, wild birds, turkey, pig) between  
443 1970 and 2016 (**Supplementary Table 1**). The total isolate collection comprised 2,425  
444 *Campylobacter* genomes, including *C. jejuni* belonging to 286 STs and 36 clonal complexes,  
445 and *C. coli* to 125 STs and 1 clonal complex. All assembled genomes and raw reads have been  
446 deposited in the NCBI repository associated with BioProject: PRJNA524315. Individual

447 accession numbers can be found in **Supplementary Table 1**. Assembled genomes of all isolates  
448 used in the study are available in FigShare DOI: 10.6084/m9.figshare.7886810.

449

#### 450 *C. coli* and *C. jejuni* phylogenies and assessing temporal signal and ‘clock-likeness’

451 Phylogenies were constructed for 601 *C. coli* and 1,824 *C. jejuni* isolates (**Supplementary**  
452 **Table 1, Figure 1**). The genomes were aligned against a reference (*C. coli* CVM N29710,  
453 accession number: NC\_022347.1 and *C. jejuni* NCTC 11168, accession number:  
454 NC\_002163.1) using MAFFT with default parameters of minimum nucleotide identity of 70%  
455 over >50% of the gene and BLAST-n word size 20. Core genes, shared by all isolates within a  
456 species (2,014 for *C. coli* and 1,668 for *C. jejuni*) were concatenated and used to construct  
457 Maximum likelihood (ML) trees using FastTree version 2.1.8 and the Generalised time-  
458 reversible (*GTR*) model of nucleotide evolution (Price et al. 2010). Isolates were analysed to  
459 test for a temporal signal of the accumulation of genetic variation over time (**Supplementary**  
460 **Figure 1**). This was carried out prior to mutation rate analysis using a phylogeny of genetic  
461 distances and sampling dates, and root to tip regression implemented in the software TempEst  
462 v1.5.1 (Rambaut et al. 2016). Core genome phylogenies contained dated-tip isolates sampled  
463 between 1970 and 2016 for *C. coli* and *C. jejuni*.

464

#### 465 *Selection of closely related isolate pairs*

466 An ideal dataset for mutation rate analysis would include isolate pairs with divergent sampling  
467 dates, sufficient to measure mutation rate over time, while remaining close enough (clustering  
468 on the tree) to share reliable recent common ancestry. Furthermore we required as many pairs  
469 as possible for confidence in average mutation rates. In order to achieve this, pairwise  
470 nucleotide identity and year between isolation date matrices were constructed separately for  
471 601 (*C. coli*) and 1,824 (*C. jejuni*) isolates. Using a bespoke R script  
472 (<https://github.com/SionBayliss/CallandMolClock>), the distribution of nucleotide identity was  
473 determined for isolate pairs within sequential isolation date categories of 1 year or more (1-37  
474 for *C. coli*, 1-46 for *C. jejuni*) by comparing every isolate to all other isolates (**Figure 2**). In  
475 each analysis, isolates were used only once as the ancestral or derived strain.

476

#### 477 *Recombination and mutation inference (quantifying nucleotide change)*

478 The raw reads of genomes (**Supplementary Table 1**) of isolate pairs (**Table 1**) were mapped  
479 to the complete reference genomes: *C. coli* YH501 (accession: NZ\_CP015528.1) and *C. jejuni*  
480 NCTC 11168 using the BWA-MEM algorithm (Li 2013). Variants were called using FreeBayes  
481 v1.1.0-dirty (Garrison and Marth 2012) and SNP effects predicted and annotated using SnpEff  
482 version 4.3 (Cingolani et al. 2012) (**Supplementary Table 4**). These tools were included in the  
483 haploid variant calling pipeline, ‘snippy’ v3.0 (<https://github.com/tseemann/snippy>). Core  
484 genome sub-tree alignments were constructed using snippy-core. Mutations introduced by point  
485 mutation and recombination were inferred on the alignments using Gubbins v2.4.1 (default  
486 settings) (Croucher et al. 2015) for each isolate pair (**Supplementary Table 4**). The snippy  
487 pipeline was used to identify synonymous and non-synonymous polymorphism within and  
488 outside of inferred recombinant regions (Croucher et al., 2015). *dN/dS* ratios were calculated  
489 for sites across the core genome using the synonymous/non-synonymous analysis program  
490 (SNAP) v2.1.1 based on the Nei and Gojobori 1986 method (Korber 2000) ([www.hiv.lanl.gov](http://www.hiv.lanl.gov)).  
491 By quantifying point mutation and recombination and synonymous and nonsynonymous  
492 polymorphism, we were able to infer different molecular evolution rate estimates. These  
493 included (i) the total mutation rate, used to calculate the number of effective lineages and (ii)  
494 the rate of accumulation of synonymous mutations occurring outside of recombinant regions,  
495 used to estimate the molecular clock. Hotspots of recombination occurring across multiple  
496 isolate pairs were observed.

497

#### 498 ***Estimating the number of coalescences at yearly intervals (Birthday problem)***

499 To consider the extent to which a given sample set represented genetic diversity within the  
500 population we developed a pipeline that calculated the number of coalescences (*effective*  
501 *lineages*, **Z**) at yearly time intervals ( $Z_1, Z_2, Z_3 \dots Z_n$ ) within the datasets. This is described by  
502 the equation  $Z = Y/X$ , Where: **Y** = all potential isolate pairs ( $n^2/2$ ); **X** = the number of possible  
503 pairs for each time interval ( $t_1, t_2, t_3 \dots t_n$ ) that is less than the predicted number of mutations  
504 that have occurred over a given time interval ( $\mu(t(1-n))$ );  $\mu$  = mutation rate;  $t$  = time interval  
505 between sampling dates, 1-46 and 1-37 years for *C. jejuni* and *C. coli* respectively. The resultant  
506 **Z** value for each time period is the estimated number of effective lineages (Birthdays) at each  
507 time cut-off, equivalent to the number of lineages sharing a common ancestor at a particular  
508 time interval (**Supplementary Figure 2**).

509

510 **Acknowledgements**

511 SKS lab funded by the Medical Research Council (MR/L015080/1, MR/S009264/1,  
512 MR/T030062/1). JKC is supported by a BBSRC-CASE studentship (BB/P504750/1). DF was  
513 supported by an MRC senior research fellowship (MR/M501608/1). EM is supported by a  
514 University of Bath Faculty of Science URSA studentship.

515

516 **Contributors**

517 JKC, SKS and DF designed the study and wrote the paper with BP. JKC, BP performed  
518 genomic analysis with input from HT, SCB and EJF. EB and MB cultured isolates for  
519 sequencing. EM, MDH and BP sequenced and assembled genomes. All authors contributed and  
520 approved the final manuscript.

521

522 **Conflict of Interest**

523 The authors declare no conflict of interest.

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542 **References**

- 543 Arnold BJ, Hanage WP. 2017. Longitudinal samples of bacterial genomes potentially bias  
544 evolutionary analyses. Unpublished data bioRxiv  
545 <https://www.biorxiv.org/content/10.1101/103465v1>, last accessed October 18, 2020  
546
- 547 Bar-On YM, Phillips R, Milo R. 2018. The biomass distribution on Earth. *Proc Natl Acad Sci*.  
548 115(25):6506–6511.
- 549
- 550 Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,  
551 Nikolenko SI, Pham S, et al. 2012. SPAdes: A new genome assembly algorithm and its  
552 applications to single-cell sequencing. *J Comput Biol*. 19(5):455-477.
- 553
- 554 Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015. Measurably evolving pathogens in the  
555 genomic era. *Trends Ecol Evol*. 30(6):306–313.
- 556
- 557 Bromham L. 2009. Why do species vary in their rate of molecular evolution? *Biol Lett*. 5:401-  
558 404.
- 559
- 560 Bronowski C, James CE, Winstanley C. 2014. Role of environmental survival in transmission  
561 of *Campylobacter jejuni*. *FEMS Microbiol Lett*. 356(1):8–19.
- 562
- 563 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Luan Wang SJL, Ruden DM. 2012. A  
564 program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:  
565 SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*. 6(2):80–92.
- 566
- 567 Cody AJ, McCarthy ND, Bray JE, Wimalarathna HML, Colles FM, Jansen van Rensburg MJ,  
568 Dingle KE, Waldenstrom J, Maiden MCJ. 2015. Wild bird-associated *Campylobacter jejuni*  
569 isolates are a consistent source of human disease, in Oxfordshire, United Kingdom. *Env Micro*  
570 *Reports* 7(5):782-788.
- 571
- 572 Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J, Van Der Linden M, Mcgee L, Von  
573 Gottberg A, Song JH, Ko KS, et al. 2011. Rapid Pneumococcal Evolution in Response to

574 Clinical Interventions. *Science* 331(6016):430–434.  
575  
576 Croucher NJ, Finkelstein JA, Pelton SI, Mitchell PK, Lee GM, Parkhill J, Bentley SD, Hanage  
577 WP, Lipsitch M. 2013. Population genomics of post-vaccine changes in pneumococcal  
578 epidemiology. *Nat Genet.* 45:656–663.  
579  
580 Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, Parkhill J, Harris SR.  
581 2015. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome  
582 sequences using Gubbins. *Nucleic Acids Res.* 43(3):e15.  
583  
584 Cui Y, Yu C, Yan Y, Li D, Li Y, Jombart T, Weinert LA, Wang Z, Guo Z, Xu L, et al. 2013.  
585 Historical variations in mutation rate in an epidemic pathogen, *Yersinia pestis*. *Proc Natl Acad*  
586 *Sci U S A.* 110(2):577–582.  
587  
588 Cui Y, Yang X, Didelot X, Guo C, Li D, Yan Y, Zhang Y, Yuan Y, Yang H, Wang J et al.  
589 2015. Epidemic clones, oceanic gene pools, and eco-LD in the free-living marine pathogen  
590 *Vibrio parahaemolyticus*. *Mol Biol Evol.* 32(6):1396-410.  
591  
592 Dearlove BL, Cody AJ, Pascoe B, Méric G, Wilson DJ, Sheppard SK. 2016. Rapid host  
593 switching in generalist *Campylobacter* strains erodes the signal for tracing human infections.  
594 *ISME J.* 10(3):721–729.  
595  
596 Denamur E, Matic I. 2006. Evolution of mutation rates in bacteria. *Mol Microbiol.* 60(4).  
597  
598 Dhillon AS, Shivaprasad HL, Schaberg D, Wier F. 2006. *Campylobacter jejuni* in broiler  
599 chickens. *Avian Diseases.* 50(1):55-58.  
600  
601 Didelot X, Falush D. 2007. Inference of Bacterial Microevolution Using Multilocus Sequence  
602 Data. *Genetics.* 175(3):1251–1266.  
603  
604 Didelot X, Eyre DW, Cule M, Ip CLC, Ansari MA, Griffiths D, Vaughan A, O’Connor L,  
605 Golubchik T, Batty EM, et al. 2012. Microevolutionary analysis of *Clostridium difficile*



606 genomes to investigate transmission. *Genome Biol.* 13(12):R118.  
607  
608 Didelot X, Nell S, Yang I, Woltemate S, Van Der Merwe S, Suerbaum S. 2013. Genomic  
609 evolution and transmission of *Helicobacter pylori* in two South African families. *Proc Natl*  
610 *Acad Sci U S A.* 110(34):13880–13885.  
611  
612 Didelot X, Walker AS, Peto TE, Crook DW, Wilson DJ. 2016. Within-host evolution of  
613 bacterial pathogens. *Nat Rev Microbiol.* 14:150-162.  
614  
615 Dingle KE, Colles FM, Wareing DRA, Ure R, Fox AJ, Bolton FE, Bootsman HJ, Willems RJL,  
616 Urwin R, Maiden MCJ. 2001. Multilocus sequence typing system for *Campylobacter jejuni*.  
617 *Journal of Clinical Microbiology* 39(1):14-23.  
618  
619 Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. 2003. Measurably evolving  
620 populations. *Trends Ecol Evol.* 18:481–488.  
621  
622 Drummond AJ, Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees.  
623 *BioMed Cent.* 7(214):1–8.  
624  
625 Duchêne S, Holmes EC, Ho SYW. 2014. Analyses of evolutionary dynamics in viruses are  
626 hindered by a time-dependent bias in rate estimates. *Proceedings Biol Sci.* 281(1786).  
627  
628 Duchêne S, Holt KE, Weill F-X, Le Hello S, Hawkey J, Edwards DJ, Fourment M, Holmes EC.  
629 2016. Genome-scale rates of evolutionary change in bacteria. *Microb Genomic.* 2(11).  
630  
631 Everitt RG, Didelot X, Batty EM, Miller RR, Knox K, Young BC, Bowden R, Auton A,  
632 Votintseva A, Lerner-Svensson H, et al. 2014. Mobile elements drive recombination hotspots  
633 in the core genome of *Staphylococcus aureus*. *Nat Commun.* 5(3956).  
634  
635 Eyre-Walker A, Keightley PD. 2007. The distribution of fitness effects of new mutations. *Nat*  
636 *Rev Gen.* 8:610-618.  
637



- 638 Flohr RCE, Blom CJ, Rainey PB, Beaumont HJE. 2013. Founder niche constrains evolutionary  
639 adaptive radiation. *PNAS*. 110(51):20663-20668.  
640
- 641 Garrison E, Marth G. 2012. Haplotype-based variant detection from short-read sequencing.  
642 Unpublished data, arXiv preprint arXiv:1207.3907 [q-bio.GN]. Last accessed on November 24  
643 2020.  
644
- 645 Gibson B, Wilson D, Feil E, Eyre-Walker A. 2018. The Distribution of Bacterial Doubling  
646 Times in the Wild. *Proc Biol Sci*. 285(1880).  
647
- 648 Gibson B, Eyre-Walker A. 2019. Investigating evolutionary rate variation in bacteria. *J Mol*  
649 *Evol*. 87:317-326.  
650
- 651 Gojobori T, Moriyama EN, Kimura M. 1990. Molecular clock of viral evolution, and the neutral  
652 theory. *Proc Natl Acad Sci U S A*. 87(24):10015-10018.  
653
- 654 Ho SYW, Shapiro B, Phillips MJ, Cooper A, Drummond AJ, Sullivan J. 2007. Evidence for  
655 Time Dependency of Molecular Rate Estimates. Sullivan J, editor. *Syst Biol*. 56(3):515–522.  
656
- 657 Ho SYW, Lanfear R, Bromham L, Phillips MJ, Soubrier J, Rodrigo AG, Cooper A. 2011. Time-  
658 dependent rates of molecular evolution. *Mol Ecol*. 20(15):3087–3101.  
659
- 660 Jolley KA, Maiden MC. 2010. BIGSdb: Scalable analysis of bacterial genome variation at the  
661 population level. *BMC Bioinformatics*. 11:595.  
662
- 663 Jolley KA, Bray JE, Maiden MCJ. 2018. Open-access bacterial population genomics: BIGSdb  
664 software, the PubMLST.org website and their applications. *Wellcome Open Res*. 24(3):124.  
665
- 666 Kidgell C, Reichard U, Wain J, Linz B, Torpdahl M, Dougan G, Achtman M. 2002. *Salmonella*  
667 *typhi*, the causative agent of typhoid fever, is approximately 50,000 years old. *Infect Genet*  
668 *Evol*. 2(1):39–45.  
669

- 670 Kimura M. 1968. Evolutionary rate at the molecular level. *Nature*. 217(5129):624-626.  
671
- 672 Kimura M. 1987. Molecular evolutionary clock and the neutral theory. *J Mol Evol*. 26(1-2):24-  
673 33.  
674
- 675 Kirchberger PC, Schmidt ML, Ochman H. 2020. The ingenuity of bacterial genomes. *Annu*  
676 *Rev Microbiol*. 8;74:815-834.  
677
- 678 Koeppl A, Perry EB, Sikorski J, Krizanc D, Warner A, Ward DM, Rooney AP, Brambilla E,  
679 Connor N, Ratcliff RM, et al. 2008. Identifying the fundamental units of bacterial diversity: A  
680 paradigm shift to incorporate ecology into bacterial systematics. *PNAS* 105(7):2504-2509.  
681
- 682 Korber B. 2000. HIV Signature and Sequence Variation Analysis. *Computational Analysis of*  
683 *HIV Molecular Sequences*, Allen G. Rodrigo and Gerald H. Learn, eds. Dordrecht,  
684 Netherlands: Kluwer Academic Publishers, Chapter 4, pages 55-72.  
685
- 686 Kuo CH, Ochman H. 2009. Inferring clocks when lacking rocks: the variable rates of molecular  
687 evolution in bacteria. *Biol Direct*. 4:35.  
688
- 689 Li H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
690 Unpublished data, [arXiv:1303.3997v2](https://arxiv.org/abs/1303.3997v2) [q-bio.GN]. Last accessed on October 18 2020  
691
- 692 Li S-J, Hua Z-S, Huang L-N, Li J, Shi S-H, Chen L-X, Kuang J-L, Liu J, Hu M, Shu W-S.  
693 2015. Microbial communities evolve faster in extreme environments. *Sci Rep*. 4(1):6205.  
694
- 695 Louca S, Shih PM, Pennell MW, Fischer WW, Parfrey LW, Doebeli M. 2018. Bacterial  
696 diversification through geological time. *Nat Ecol Evol*. 2, 1458-1467.  
697
- 698 Maixner F, Krause-Kyora B, Turaev D, Herbig A, Hoopmann MR, Hallows JL, Kusebauch U,  
699 Egarter Vigl E, Malfertheiner P, Megraud F, et al. 2016. The 5,300-year-old *Helicobacter pylori*  
700 genome of the Iceman HHS Public Access. *Science* 351(6269):162-165.  
701

- 702 Marin J, Hedges SB. 2018. Undersampling genomes has biased time and rate estimates  
703 throughout the tree of life. *Mol Biol Evol.* 35(10):2595.  
704
- 705 Mathers AJ, Stoesser N, Sheppard AE, Pankhurst L, Giess A, Yeh AJ, Didelot X, Turner SD,  
706 Sebra R, Kasarskis A, et al. 2015. *Klebsiella pneumoniae* Carbapenemase (KPC)-Producing *K.*  
707 *pneumoniae* at a Single Institution: Insights into Endemicity from Whole-Genome Sequencing.  
708 *Antimicrob Agents Chemother.* 59:1656–1663.  
709
- 710 Mathis, FH. 1991. A generalized birthday problem. *SIAM review.* 33(2):265-270.  
711
- 712 Mcadam PR, Templeton KE, Edwards GF, Holden MTG, Feil EJ, Aanensen DM, Bargawi  
713 HJA, Spratt BG, Bentley SD, Parkhill J, et al. 2012. Molecular tracing of the emergence,  
714 adaptation, and transmission of hospital-associated methicillin-resistant *Staphylococcus*  
715 *aureus*. *Proc Natl Acad Sci U S A.* 109(23):9107–9112.  
716
- 717 Menardo F, Duchêne S, Brites D, Gagneux S. 2019. The molecular clock of *Mycobacterium*  
718 *tuberculosis*. *PLoS Pathog.* 15(9):e1008067.  
719
- 720 von Mering C, Hugenholtz P, Raes J, Tringe SG, Doerks T, Jensen LJ, Ward N, Bork P. 2007.  
721 Quantitative Phylogenetic Assessment of Microbial Communities in Diverse Environments.  
722 *Science* 315(5815):1126–1130.  
723
- 724 Mourkas E, Taylor AJ, Meric G, Bayliss SC, Pascoe B, Mageiros L, Calland JK, Hitchings MD,  
725 Ridley A, Vidal A, et al. 2020. Agricultural intensification and the evolution of host specialism  
726 in the enteric pathogen *Campylobacter jejuni*. *PNAS* 117(20):11018-11028.  
727
- 728 Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, Croucher NJ, Choi SY,  
729 Harris SR, Lebens M, et al. 2011. Evidence for several waves of global transmission in the  
730 seventh cholera pandemic. *Nature* 477(7365):462–465.  
731
- 732 Nei M, Gojobori T. 1986. Simple methods for estimating the numbers of synonymous and  
733 nonsynonymous nucleotide substitutions. *Mol Biol Evol.* 3(5):418-426.

734

735 Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees  
736 for Large Alignments. Poon AFY, editor. PLoS One. 5(3):e9490.

737

738 Rainey PB, Travisano M. 1998. Adaptive radiation in a heterogeneous environment. *Nature*.  
739 394(6688):69-72.

740

741 Rambaut A, Lam TT, Carvalho LM, Pybus OG. 2016. Exploring the temporal structure of  
742 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* 2(1).

743

744 Rieux A, Balloux F. 2016. Inferences from tip-calibrated phylogenies: a review and a practical  
745 guide. *Mol Ecol.* 25(9):1911–1924.

746

747 Rocha EPC, Smith JM, Hurst LD, Holden MTG, Cooper JE, Smith NH, Feil EJ. 2006.  
748 Comparisons of dN/dS are time dependent for closely related bacterial genomes. *J Theor Biol.*  
749 239:226–235.

750

751 Rosner BM, Schielke A, Didelot X, Kops F, Breidenbach J, Willrich N, Götz G, Alter T, Stingl  
752 K, Josenhans C, et al. 2017. A combined case-control and molecular source attribution study of  
753 human *Campylobacter* infections in Germany, 2011–2014. *Sci Rep.* 7(1):5139.

754

755 Sheppard SK, McCarthy ND, Falush D, Maiden MCJ. 2008. Convergence of *Campylobacter*  
756 species: implications for bacterial evolution. *Science* 320(5873):237-239.

757

758 Sheppard SK, Dallas JF, Strachan NJC, MacRae M, McCarthy ND, Wilson DJ, Gormley FJ,  
759 Falush D, Ogden ID, Maiden MCJ, et al. 2009. *Campylobacter* genotyping to determine the  
760 source of human infection. *Clin Infect Dis.* 48(8):1072–1078.

761

762 Sheppard SK, Dallas JF, Wilson DJ, Strachan NJC, McCarthy ND, Jolley KA, Colles FM,  
763 Rotariu O, Ogden ID, Forbes KJ, et al. 2010. Evolution of an Agriculture-Associated Disease  
764 Causing *Campylobacter coli* Clade: Evidence from National Surveillance Data in Scotland.  
765 Hartskeerl RA, editor. PLoS One. 5(12):e15708.

766

767 Sheppard SK, Colles FM, Mccarthy ND, Strachan NJC, Ogden ID, Forbes KJ, Dallas JF,  
768 Maiden MCJ. 2011. Niche segregation and genetic structure of *Campylobacter jejuni*  
769 populations from wild and agricultural host species. *Mol Ecol.* 20(16):3484–3490.

770

771 Sheppard SK, Didelot X, Jolley KA, Darling AE, Pascoe B, Meric G, Kelly DJ, Cody A, Colles  
772 FM, Strachan NJC, et al. 2013. Progressive genome-wide introgression in agricultural  
773 *Campylobacter coli*. *Mol Ecol.* 22(4):1051-1064.

774

775 Sheppard SK, Cheng L, Méric G, De Haan CPA, Llarena AK, Marttinen P, Vidal A, Ridley A,  
776 Clifton-Hadley F, Connor TR, et al. 2014. Cryptic ecology among host generalist  
777 *Campylobacter jejuni* in domestic animals. *Mol Ecol.* 23(10):2442–2451.

778

779 Sheppard SK, Maiden MCJ. 2015. The evolution of *Campylobacter jejuni* and *Campylobacter*  
780 *coli*. *Cold Spring Harb Perspec Biol* 7(8):a018119.

781

782 Strachan NJC, Gormley FJ, Rotariu O, Ogden ID, Miller G, Dunn GM, Sheppard SK, Dallas  
783 JF, Reid TMS, Howie H, et al. 2009. Attribution of *Campylobacter* Infections in Northeast  
784 Scotland to Specific Sources by Use of Multilocus Sequence Typing. *J Infect Dis.* 199(8):1205–  
785 1208.

786

787 Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018. Bayesian  
788 phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol* 8:4(1):vey016.

789

790 Thakur S, Morrow WEM, Funk JA, Bahnson PB, Gebreyes WA. 2006. Molecular  
791 epidemiologic investigation of *Campylobacter coli* in swine production systems, using  
792 multilocus sequence typing. *Appl Environ Microbiol.* 72(8):5666-5669.

793

794 Thépault A, Méric G, Rivoal K, Pascoe B, Mageiros L, Touzain F, Rose V, Béven V, Chemaly  
795 M, Sheppard SK. 2017. Genome-Wide Identification of Host- Segregating Epidemiological  
796 Markers for Source Attribution in *Campylobacter jejuni*. *Appl Environ Microbiol.* 83(7).

797

798 Vos M, Didelot X. 2009. A comparison of homologous recombination rates in bacteria and  
799 archaea. *ISME J.* 3(2):199–208.  
800

801 Waldenström J, Broman T, Carlsson I, Hasselquist D, Achterberg RP, Wagenaar JA, Olsen B.  
802 2002. Prevalence of *Campylobacter jejuni*, *Campylobacter lari*, and *Campylobacter coli* in  
803 different ecological guilds and taxa of migrating birds. *Appl Environ Microbiol.* 68(12):5911–  
804 5917.  
805

806 Walker TM, C Ip CL, Harrell RH, Evans JT, Kapatai G, Dedicoat MJ, Eyre DW, Wilson DJ,  
807 Hawkey PM, Crook DW, et al. 2013. Whole-genome sequencing to delineate *Mycobacterium*  
808 tuberculosis outbreaks: a retrospective observational study. *Lancet Infect Dis.* 13:137–146.  
809

810 Weller C, Wu M. 2015. A generation-time effect on the rate of molecular evolution in bacteria.  
811 *Evolution (N Y).* 69(3):643–652.  
812

813 Wilson DJ, Gabriel E, Leatherbarrow AJH, Cheesbrough J, Gee S, Bolton E, Fox A, Fearnhead  
814 P, Hart CA, Diggle PJ. 2008. Tracing the Source of Campylobacteriosis. *PLoS Genet.*  
815 4(9):e1000203.  
816

817 Wilson DJ, Gabriel E, Leatherbarrow AJH, Cheesbrough J, Gee S, Bolton E, Fox A, Hart CA,  
818 Diggle PJ, Fearnhead P. 2009. Rapid evolution and the importance of recombination to the  
819 gastroenteric pathogen *Campylobacter jejuni*. *Mol Biol Evol.* 26(2):385–397.  
820

821 Yahara K, Didelot X, Ansari MA, Sheppard SK, Falush D. 2014. Efficient inference of  
822 recombination hot regions in bacterial genomes. *Molecular Biology and Evolution* 31(6): 1593-  
823 1605.  
824  
825  
826  
827  
828  
829

830 **Figure legends**

831

832 **Figure 1. Little evidence of clustering of isolate sampling dates in *Campylobacter***  
833 **phylogenies.** Maximum likelihood (ML) core genome phylogenetic trees of *C. coli* (**A**) (n =  
834 601) and *C. jejuni* (**B**) (n = 1824) constructed using FastTree version 2.1.8 (Price et al. 2010)  
835 and the *GTR* model of nucleotide evolution. Both phylogenies show the distribution of the  
836 sample time frame used in this study with major *Campylobacter* clonal complexes identified  
837 and terminal nodes coloured according to isolation decade (orange = 1970s, yellow = 1980s,  
838 white = 1990s, green = 2000s, blue = 2010s). Scale bars represent the estimated number of  
839 mutations per site. Terminal nodes sampled from different decades can be seen scattered  
840 throughout both trees with little evidence of clustering by decade. Isolates sampled from the  
841 2000s and 2010s are most abundant within each dataset.

842

843 **Figure 2. Pair selection criteria curves for inclusion in rate estimates.** Visual representation  
844 of possible pairs of isolates at all time cut-offs across the sample time frame for *C. jejuni* (**B**)  
845 and *C. coli* (**D**). As time difference between pairs increases, distinguishing between individual  
846 curves becomes distorted. Therefore, a selection of years were plotted (**A** and **C**) (black = all  
847 pairs >1 year difference, pink = >2 years, blue = >4 years, purple = >6 years, orange = >8 years,  
848 red = >10 years, green = >15 years). All isolates were paired with the nearest isolate (genetic  
849 distance), matched according to difference in year of isolation (coloured lines) for both *C. jejuni*  
850 (**A**) and *C. coli* (**C**) (orange line). Dashed boxes (**A** and **C**) show magnified images of the closest  
851 pairs from all curves. Grey scale bars (**B** and **D**) indicate the time difference cut-off of each  
852 curve for every time point in the sample date frame.

853

854 **Figure 3. Mutation and recombination in *C. coli* and *C. jejuni*.** Average genome-wide SNP  
855 positions (red dots = synonymous polymorphisms, blue dots = nonsynonymous  
856 polymorphisms) per isolate pair in relation to inferred recombined regions (grey blocks). Each  
857 plot represents one pair of isolates considered in rate calibration for *C. coli* (**A**) and *C. jejuni*  
858 (**B**) and are ordered according to **Table 1**. *y axis* = number of substitutions in relation to  
859 particular bp position of the reference genome (*C. coli* = YH501, *C. jejuni* = NCTC11168) and  
860 varies between pairs. *x axis* = position of reference genome in bins of 10,000 bp. The cladogram  
861 shows the relatedness of isolate pairs based on nucleotide identity, scale bar indicates



862 polymorphisms per site. It is evident from both **A** and **B** that recombination is the main source  
863 of variation in *C. coli* and *C. jejuni*.

864

865 **Figure 4. Lineage expansion in *C. jejuni* and *C. coli*.** (A) Number of effective lineages (y  
866 axis) at each time point within the sample time frame (x axis) for *C. coli* (grey) and *C. jejuni*  
867 (black). (B) Diagrammatic representation of lineage expansion in *C. coli* and *C. jejuni* showing  
868 contrasting lineage diversification scenarios.

869

870 **Supplementary Figure 1.** Root-to-tip linear regression of *C. coli* and *C. jejuni* implemented in  
871 the software, TempEst. Root-to-tip genetic distance (y axis) is correlated against sampling times  
872 (x axis) for phylogenies of 601 *C. coli* (A) and 1,824 *C. jejuni* (B). Although both *C. coli* and  
873 *C. jejuni* datasets show a weak temporal signal, positive correlations can be seen for both  
874 species.

875

876 **Supplementary Figure 2.** Methods for calculating the number of effective lineages within the  
877 population. (A) The total number of *C. jejuni* and *C. coli* isolates in the population and all  
878 potential pairwise comparisons between putative ancestral (black) and contemporary (white)  
879 strains to give the total number of potential isolate pairs, **Y**. (B) Isolate pair selection based on  
880 divergent sampling date (>8 years) and a nucleotide identity threshold <5000 SNPs. (C) Total  
881 mutation rate ( $\mu$ ) calculated for all chosen pairs. The rate of accumulation of all synonymous  
882 (Sd), nonsynonymous (Sn) substitutions, within (rec) and outside (mut) of recombined regions,  
883 was estimated since the most recent common ancestor (MRCA, red circle). The difference in  
884 substitutions between each pair was divided by the difference in isolation years to give  $\mu$ . (D)  
885 The mutation rate was used to estimate the number of SNPs that were to accumulate over a time  
886 period and the number of possible isolate pairs at given time intervals ( $t_1, t_2, t_3, \dots, t_n$ ) for each  
887 species.

888

889 **Table titles**

890

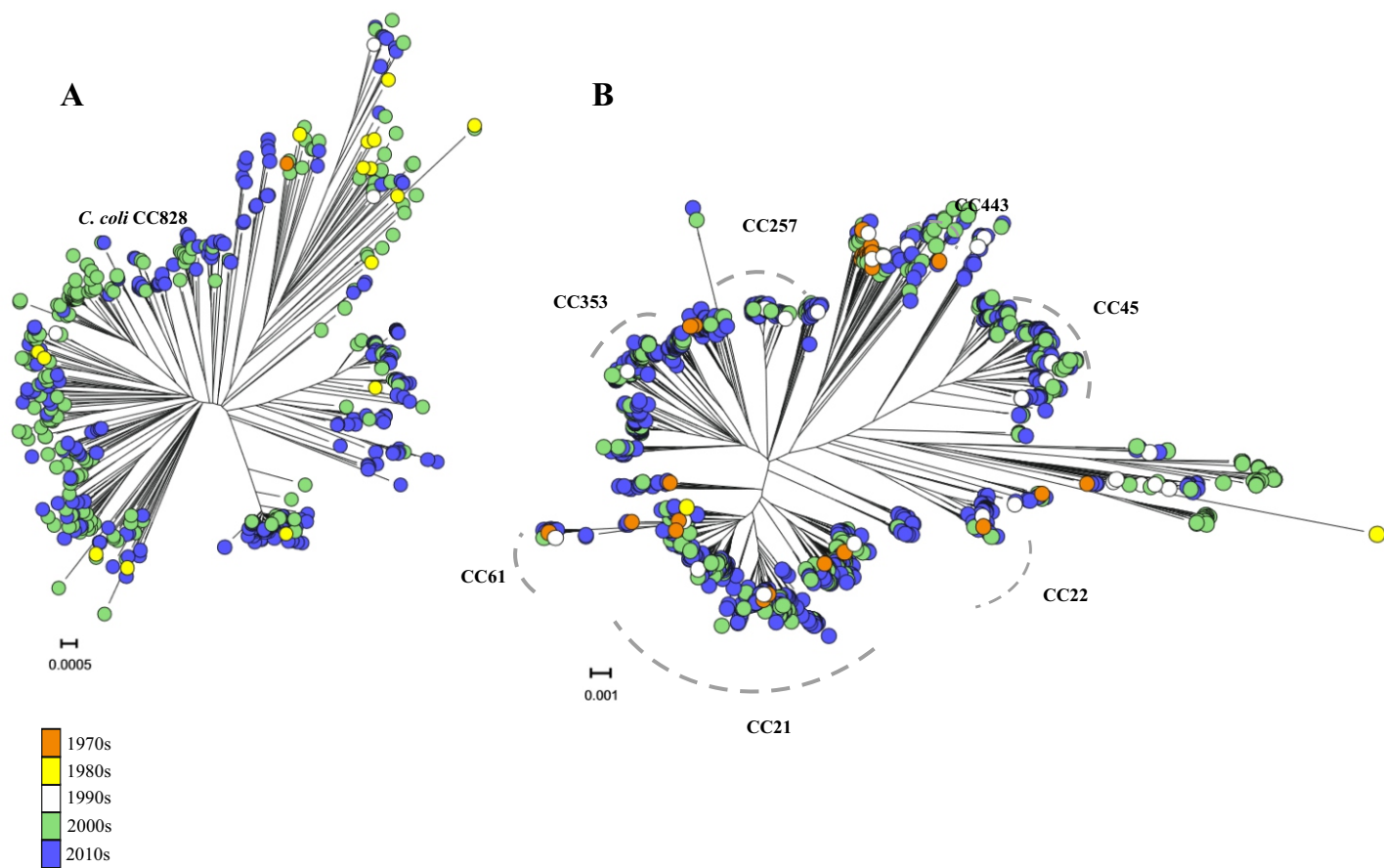
891 **Table 1: Isolate pair information for *C. coli***

892 **Table 2: Isolate pair information for *C. jejuni***



- 893 **Table 3: Estimates of evolutionary potential of nucleotide change across all *C. coli***  
894 **genomes in each pair of isolates**
- 895 **Table 4: Estimates of evolutionary potential of nucleotide change across all *C. jejuni***  
896 **genomes in each pair of isolates**
- 897 **Table 5: Average rate calibrations in *C. jejuni* and *C. coli***
- 898 **Table 6: Recombination information for each isolate in each *C. coli* pair considered for**  
899 **rate calibration**
- 900 **Table 7: Recombination information for each isolate in each *C. jejuni* pair considered**  
901 **for rate calibration**
- 902 **Table 8: *C. coli* “birthday problem” data and estimates of coalescence across sample**  
903 **time frame**
- 904 **Table 9: *C. jejuni* “birthday problem” data and estimates of coalescence across sample**  
905 **time frame**
- 906
- 907 **Supplementary Table 1: Isolate list information**
- 908 **Supplementary Table 2: TempEst root-tip regression analysis estimates**
- 909 **Supplementary Table 3: List of possible pairs of isolates**
- 910 **Supplementary Table 4: Additional SNP annotations**
- 911 **Supplementary Table 5: Individual pair rates**

Figure 1.



## Figure 2.

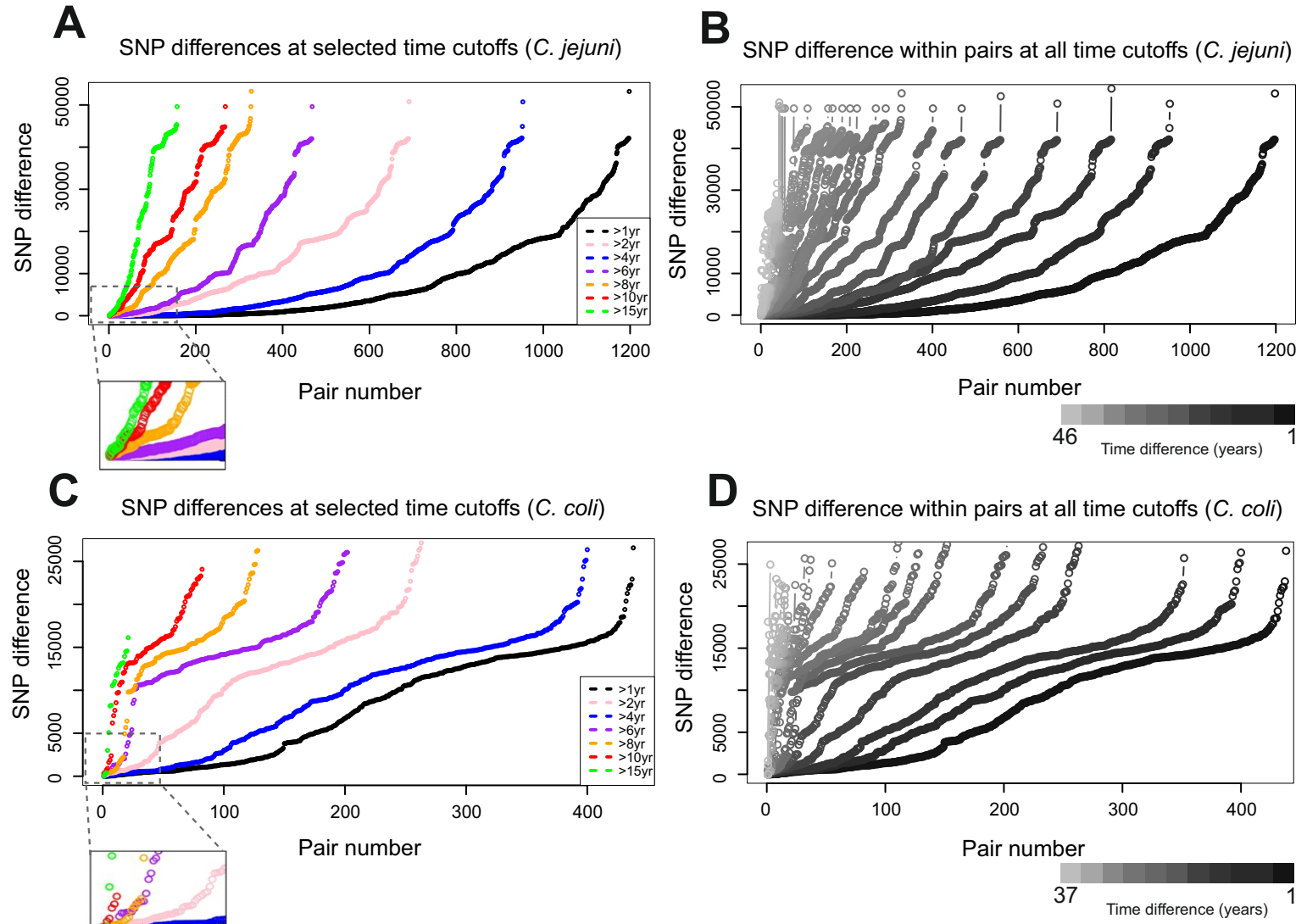


Figure 3.

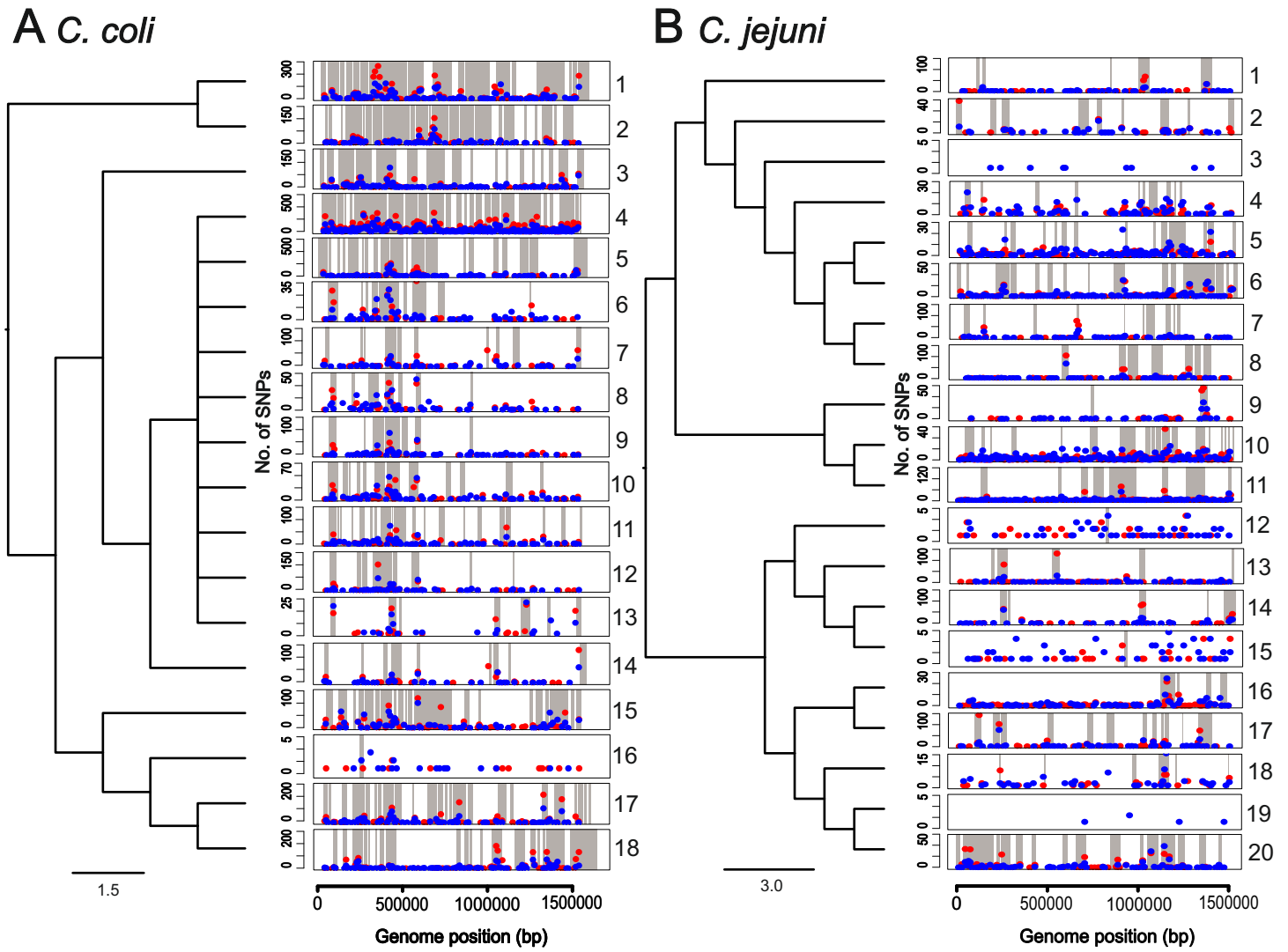
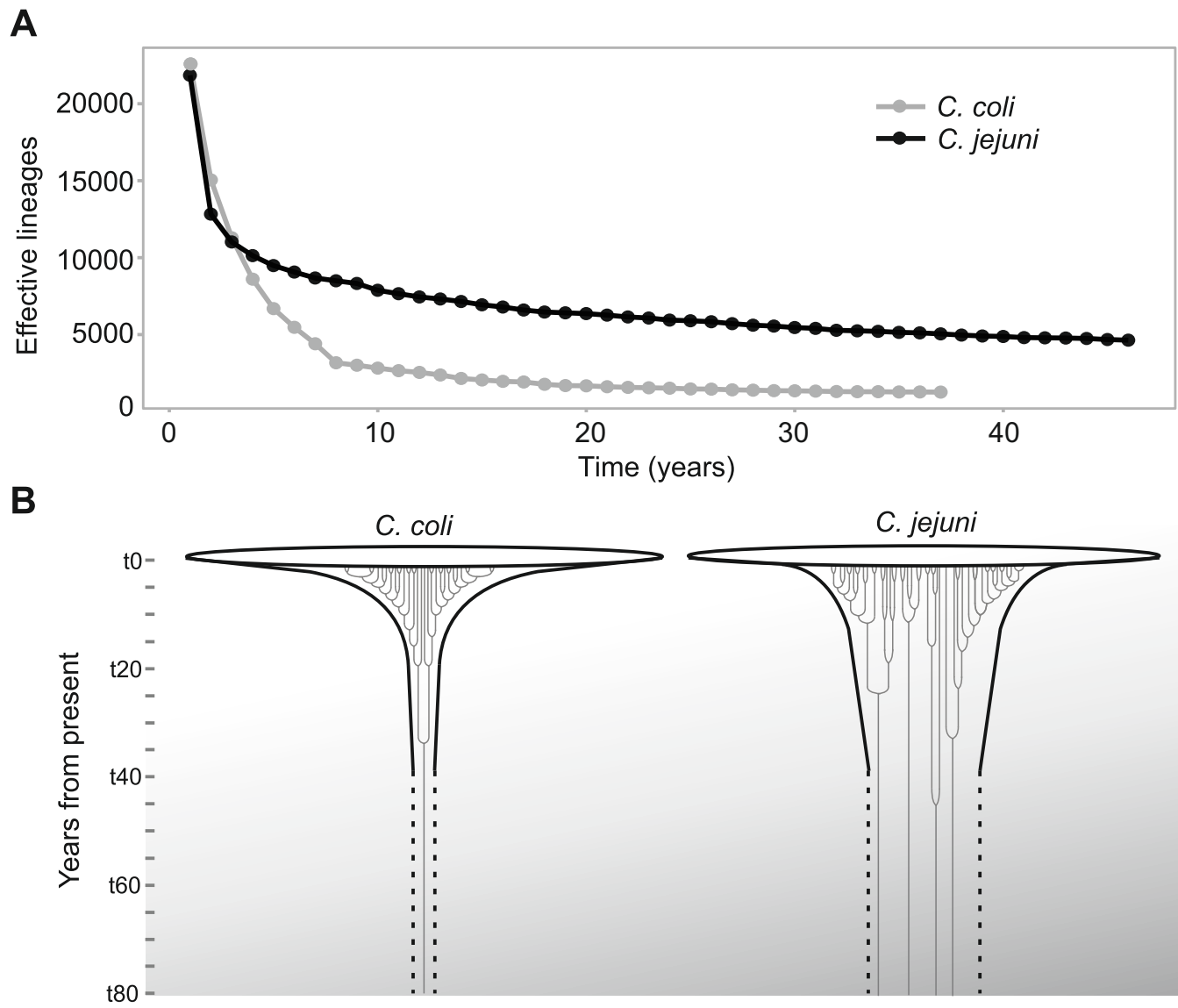
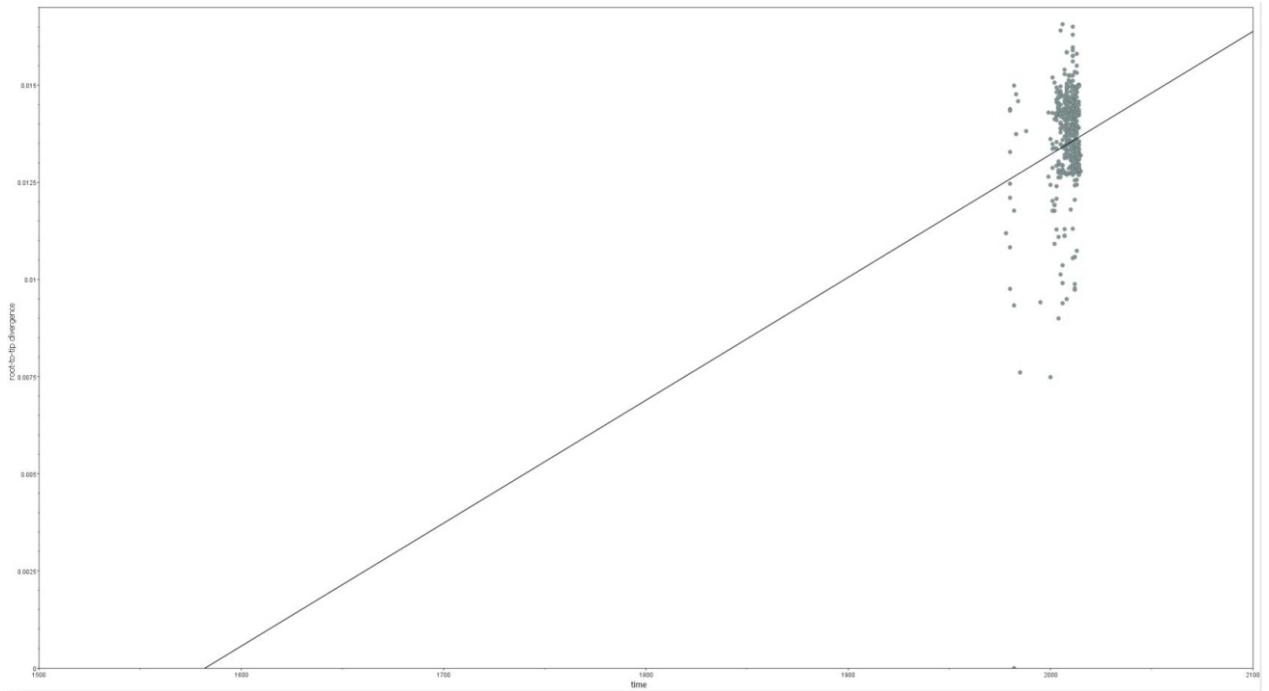


Figure 4.

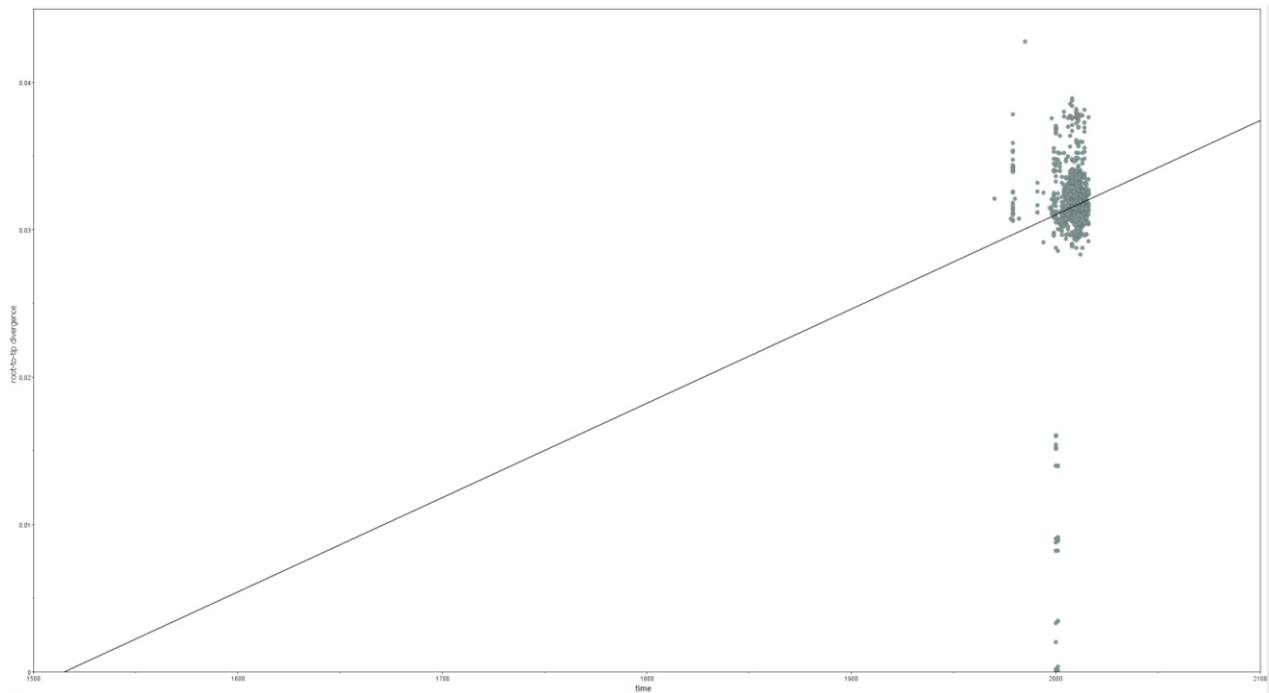


# Supplementary Figure 1.

## A) *C. coli*

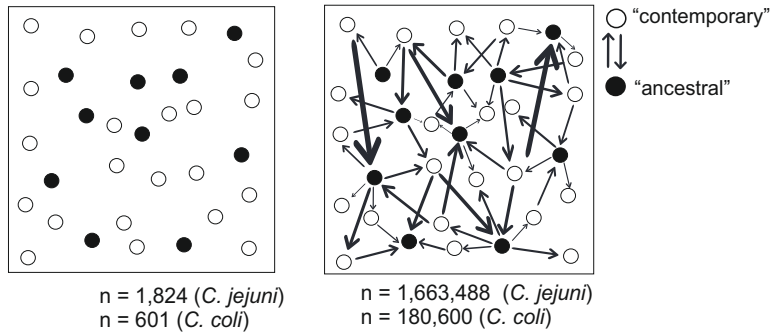


## B) *C. jejuni*

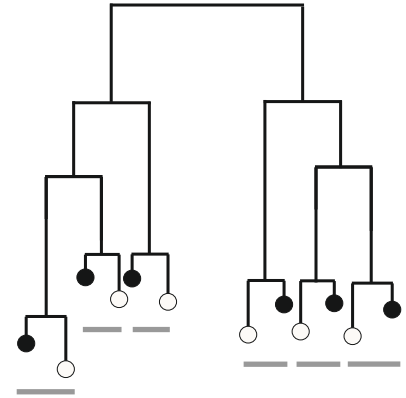


## Supplementary Figure 2.

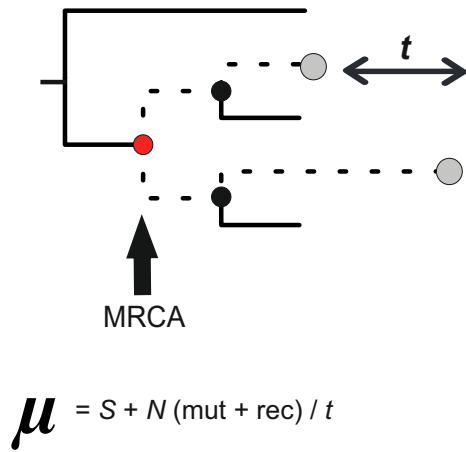
A)



B)



C)



D)

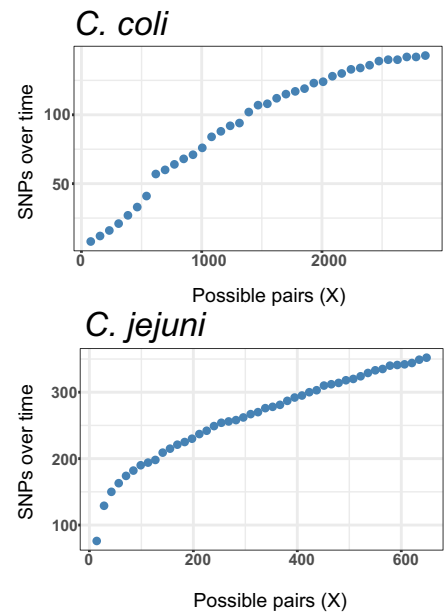


Table 1: Isolate pair information for *Campylobacter coli*

Pair no.	Isolate	Year of isolation	ST*	CC**	Source	Country	Difference in years
1	4316.LDI12946	2006	827	ST-828CC	Environmental	UK	8
	3158.LDI6744	2014	827	ST-828CC	Environmental	UK	
2	4281.LDI12911	2007	827	ST-828CC	Environmental	UK	8
	3804.CCN182coli	2015	827	ST-828CC	Duck	UK	
3	1783.SS_328	2006	827	ST-828CC	Cattle	UK	8
	3166.LDI6752	2014	827	ST-828CC	Environmental	UK	
4	454.SS_047	2005	827	ST-828CC	Chicken	UK	8
	3150.LDI6736	2013	827	ST-828CC	Environmental	UK	
5	3899.H042120298	2004	962	ST-828CC	Environmental	UK	8
	2701.OXC6817	2012	962	ST-828CC	Human	UK	
6	1770.SS_018	2006	827	ST-828CC	Chicken	UK	8
	3160.LDI6746	2014	827	ST-828CC	Environmental	UK	
7	439.SS_031	2005	827	ST-828CC	Chicken	UK	8
	3152.LDI6738	2013	827	ST-828CC	Environmental	UK	
8	1789.SS_335	2006	827	ST-828CC	Wild bird	UK	8
	3174.LDI6760	2014	827	ST-828CC	Environmental	UK	
9	1798.SS_344	2001	827	ST-828CC	Chicken	UK	10
	1753.SS_273	2011	827	ST-828CC	Human	UK	
10	3766.H065100499	2006	827	ST-828CC	Environmental	UK	8
	3176.LDI6762	2014	827	ST-828CC	Environmental	UK	
11	3922.H054900335	2005	825	ST-828CC	Environmental	UK	8
	3161.LDI6747	2013	825	ST-828CC	Environmental	UK	
12	4348.LDI12978	2005	827	ST-828CC	Environmental	UK	8
	3154.LDI6740	2013	827	ST-828CC	Environmental	UK	
13	3925.H043900429	2004	827	ST-828CC	Environmental	UK	8
	2690.OXC6785	2012	827	ST-828CC	Human	UK	
14	3877.H054000445	2005	825	ST-828CC	Environmental	UK	8
	3169.LDI6755	2013	825	ST-828CC	Environmental	UK	
15	1797.SS_343	2005	825	ST-828CC	Chicken	UK	8
	3162.LDI6748	2013	825	ST-828CC	Environmental	UK	
16	3800.UNOR5482c	2003	2195	ST-828CC	Environmental	UK	9
	1855.SS_614	2012	2195	ST-828CC	Human	UK	
17	3803.UNOR13691b	2003	827	ST-828CC	Environmental	UK	8
	1792.SS_338	2011	827	ST-828CC	Chicken	UK	
18	3167.LDI6753	2001	1541	ST-828CC	Environmental	UK	11
	1852.SS_595	2012	6795	ST-828CC	Human	UK	



Table 2: Isolate pair information for *Campylobacter jejuni*

Pair no.	Isolate	Year of isolation	ST*	CC**	Source	Country	Difference in years
1	5932.Manchester	2003	43	ST-21CC	Human	UK	13
	29.NC_002163	2016	43	ST-21CC	Human	UK	
2	275.13264	1999	257	ST-257CC	Human	UK	12
	2140.OXC5779	2011	257	ST-257CC	Human	UK	
3	267.13256	1991	42	ST-42CC	Human	UK	12
	85.cow3583	2003	3583	ST-42CC	Cattle	UK	
4	5921.Aberdeen	2002	43	ST-21CC	Human	UK	14
	5855.LITTER.B.E	2016	43	ST-21CC	Chicken	UK	
5	269.13258	1998	48	ST-48CC	Lamb	UK	9
	1806.SS_381	2007	48	ST-48CC	Chicken	UK	
6	1776.SS_029	2005	45	ST-45CC	Chicken	UK	10
	4809.CTA244	2015	45	ST-45CC	Dog	France	
7	5735.79.23	1979	1457	ST-443CC	Duck	USA	32
	2397.Seal73	2011	1457	ST-443CC	Seal	UK	
8	255.Hn30	2003	50	ST-21CC	Human	UK	12
	4801.CTA093	2015	50	ST-21CC	Dog	France	
9	299.SS_202	2008	257	ST-257CC	Chicken	UK	8
	5814.A4.G	2016	257	ST-257CC	Chicken	UK	
10	303.SS_214	2008	48	ST-48CC	Chicken	UK	8
	5986.LE.72	2016	48	ST-48CC	Chicken	UK	
11	268.13257	1999	45	ST-45CC	Human	UK	12
	2533.OXC6314	2011	45	ST-45CC	Human	UK	
12	5592.RM1285	1997	50	ST-21CC	Chicken	USA	16
	5560.MTVDSCj07	2013	50	ST-21CC	Chicken	USA	
13	2928.SS_0890	1999	1326	ST-45CC	Wild Bird	Sweden	12
	2451.Seal186	2011	1326	ST-45CC	Seal	UK	
14	5038.CjRM3147	2001	22	ST-22CC	Human	Mexico	15
	5902.Camp_108c	2016	22	ST-22CC	Chicken	UK	
15	1800.SS_375	2007	45	ST-45CC	Chicken	UK	8
	4811.CTA277	2015	45	ST-45CC	Dog	France	
16	463.SS_058	2005	45	ST-45CC	Chicken	UK	10
	5006.SS_2784	2015	1701	ST-45CC	Chicken	UK	
17	62.cowa21	2006	21	ST-21CC	Cattle	UK	9
	4834.CTA710	2015	21	ST-21CC	Dog	France	
18	5738.79.248	1979	50	ST-21CC	Turkey	USA	30
	5595.WP2202	2009	50	ST-21CC	Chicken	USA	
19	272.13261	1998	61	ST-61CC	Cattle	UK	18
	191.Cj2008.872	2016	61	ST-61CC	Human	France	
20	5743.79.315	1979	50	ST-21CC	Human	USA	36
	4832.CTA693	2015	50	ST-21CC	Dog	France	

Table 3: Estimates of evolutionary potential of nucleotide change across all *C. coli* genomes in each pair of isolates

Pair no.	SNPs in rec		SNPs out rec		rdN/dS	mdN/dS	SNPs/year*	mol. clock rates**	Total SNPs rate (m + r)/yr***
	S	N	S	N					
1	202	249	61	105	0.731	0.541	-17.375	-6.000	-78.375
	0	0	13	20	0.000	0.575			
2	0	0	0	0	0.000	0.595	6.500	1.875	32.000
	104	95	15	34	0.473	0.597			
3	167	170	8	14	0.474	0.564	6.000	2.000	45.500
	276	349	24	31	0.362	0.558			
4	0	0	0	0	0.000	0.540	23.125	6.250	65.000
	171	157	50	113	0.540	0.619			
5	0	0	0	0	0.000	0.595	20.000	10.250	171.750
	609	569	82	64	0.564	0.509			
6	110	138	17	24	0.578	0.653	-2.125	-1.375	28.375
	199	284	6	17	0.376	0.571			
7	74	96	8	10	0.744	0.542	5.125	1.375	54.875
	350	210	19	35	0.692	0.579			
8	283	295	14	16	0.938	0.611	2.875	0.375	3.125
	247	314	17	23	0.288	0.543			
9	11717	6313	2467	1520	0.598	0.528	-256.400	-151.200	-1228.200
	4905	3639	955	644	1.144	1.115			
10	40	44	1	0	0.482	0.554	4.000	1.125	69.750
	282	316	10	24	0.679	0.611			
11	7	7	11	10	1.177	0.566	-1.875	-1.250	-3.625
	2	1	1	4	0.000	0.585			
12	0	0	0	0	0.000	0.620	10.250	4.125	62.250
	257	140	33	45	0.527	0.575			
13	49	75	36	74	0.889	0.622	10.875	5.250	177.250
	702	645	78	112	0.575	0.562			
14	419	262	19	44	0.689	0.750	10.000	3.125	107.125
	899	524	44	90	0.720	0.548			
15	947	622	79	39	0.458	0.613	-8.875	-8.125	-174.250
	191	107	14	37	0.354	0.500			
16	549	370	58	64	0.493	0.581	-15.333	-6.444	-126.111
	0	0	0	0	0.464	0.588			
17	681	510	60	129	0.611	0.618	-3.125	1.625	110.500
	1102	932	73	96	0.592	0.580			
18	0	0	0	0	0.000	0.592	27.273	13.818	545.909
	3304	2220	152	127	0.515	0.474			
<b>SD</b>	2101.48	1204.22	431.68	266.48	0.311	0.101	62.649	36.435	339.754
<b>Mean</b>	801	546	123	99	0.492	0.594	11.457	4.266	113.339

\*synonymous and nonsynonymous substitution outside of recombination; \*\*synonymous substitutions outside of recombination; \*\*\*all substitutions within and without of recombination

S = observed synonymous mutations, N = observed nonsynonymous mutations, SD = standard deviation

Table 4: Estimates of evolutionary potential of nucleotide change across all *C. jejuni* genomes in each pair of isolates

Pair no.	SNPs in rec		SNPs out rec		rdN/dS	mdN/dS	SNPs/year*	mol. clock rates**	Total SNPs rate (m + r)/yr***
	S	N	S	N					
1	0	0	0	3	0.000	0.000	0.000	0.000	0.077
	0	1	0	1	0.000	0.000			
2	30	50	23	40	0.488	0.505	0.167	0.083	-6.750
	0	0	24	47	0.000	0.580			
3	0	0	0	0	0.000	0.521	9.667	3.250	10.667
	4	3	39	67	0.342	0.517			
4	0	0	0	1	0.000	0.000	0.571	0.000	0.571
	0	0	0	8	0.000	0.000			
5	314	157	33	48	0.462	0.565	-8.000	-2.556	-57.444
	18	20	10	14	2.236	0.524			
6	136	205	90	174	0.730	0.545	58.300	26.800	32.100
	20	62	358	432	0.649	0.511			
7	0	0	0	0	0.000	0.524	7.219	3.000	10.844
	52	61	96	118	0.667	0.534			
8	16	47	32	32	1.131	0.544	27.583	10.917	39.583
	67	128	163	212	0.684	0.569			
9	107	63	20	36	0.000	0.563	-5.500	-1.875	-27.250
	1	2	5	6	1.225	0.503			
10	0	0	1	4	0.000	0.500	7.125	1.875	27.125
	97	62	16	40	0.459	0.534			
11	0	0	0	0	0.000	0.534	9.583	2.750	10.250
	0	7	33	77	0.786	0.542			
12	242	114	50	61	0.635	0.547	-6.313	-2.875	-28.000
	3	8	4	11	0.605	0.598			
13	259	138	27	42	0.311	0.560	-5.500	-2.083	-39.667
	0	0	2	8	0.000	0.538			
14	31	29	2	11	0.406	0.515	4.333	2.067	13.067
	128	61	33	40	0.639	0.571			
15	370	211	64	88	0.508	0.513	45.125	23.500	-20.000
	16	50	252	288	0.322	0.577			
16	172	69	79	101	0.399	0.531	-19.600	-7.500	-44.400
	0	0	4	9	0.000	0.500			
17	37	62	117	210	0.550	0.565	-39.222	-13.000	-51.000
	0	0	0	0	0.000	0.521			
18	10	26	0	0	1.265	0.578	2.300	0.367	10.267
	174	88	11	47	0.632	0.542			
19	177	128	22	44	0.840	0.511	-0.833	0.056	0.556
	202	119	23	33	0.441	0.520			
20	270	201	24	39	0.639	0.482	3.972	1.278	-0.167
	121	204	70	125	0.576	0.566			
<b>SD</b>	102.78	67.13	72.07	89.59	0.461	0.165	20.772	9.085	27.767
<b>Mean</b>	77	59	43	63	0.490	0.509	13.534	5.424	14.101

\*synonymous and nonsynonymous substitution outside of recombination; \*\*synonymous substitutions outside of recombination; \*\*\*all substitutions within and without of recombination

S = observed synonymous mutations, N = observed nonsynonymous mutations, SD = standard deviation

**Table 5: Average rate calibrations in *C. coli* and *C. jejuni***

	Units**	<i>C. coli</i>			<i>C. jejuni</i>		
		Min	Mean	Max	Min	Mean	Max
<b>Total substitution rate*</b>	s/s/y	1.7 x 10 <sup>-6</sup>	6.3 x 10 <sup>-5</sup>	3.0 x 10 <sup>-4</sup>	4.8 x 10 <sup>-8</sup>	8.8 x 10 <sup>-6</sup>	2.3 x 10 <sup>-5</sup>
<b>Total substitution rate absent of recombining sequences</b>	s/s/y	1.6 x 10 <sup>-6</sup>	6.4 x 10 <sup>-6</sup>	1.5 x 10 <sup>-6</sup>	1.0 x 10 <sup>-7</sup>	8.5 x 10 <sup>-6</sup>	3.6 x 10 <sup>-5</sup>
<b>Synonymous substitution rate in recombining sequences</b>	s/s/y	4.9 x 10 <sup>-6</sup>	3.1 x 10 <sup>-4</sup>	1.8 x 10 <sup>-3</sup>	2.1 x 10 <sup>-7</sup>	1.9 x 10 <sup>-6</sup>	7.6 x 10 <sup>-6</sup>
<b>Synonymous mutation rate absent of recombining sequences (molecular clock)</b>	s/s/y	2.1 x 10 <sup>-7</sup>	2.4 x 10 <sup>-6</sup>	7.7 x 10 <sup>-6</sup>	3.8 x 10 <sup>-8</sup>	3.4 x 10 <sup>-6</sup>	1.7 x 10 <sup>-5</sup>
<b>Nonsynonymous substitution rate in recombining sequences</b>	s/s/y	4.4 x 10 <sup>-8</sup>	2.4 x 10 <sup>-5</sup>	1.1 x 10 <sup>-4</sup>	5.0 x 10 <sup>-8</sup>	1.4 x 10 <sup>-6</sup>	4.8 x 10 <sup>-6</sup>
<b>Nonsynonymous mutation rate absent of recombining sequences</b>	s/s/y	4.8 x 10 <sup>-7</sup>	3.2 x 10 <sup>-6</sup>	7.8 x 10 <sup>-6</sup>	3.1 x 10 <sup>-7</sup>	4.8 x 10 <sup>-6</sup>	1.6 x 10 <sup>-5</sup>

\*all substitutions from within and outside recombination

\*\*substitutions per site per year (*C. jejuni* = 1.6 Mbp, *C. coli* = 1.8 Mbp)

Table 6: Recombination information for each isolate in each *C. coli* pair considered for rate calibration

Pair no.	Rec blocks from root of subtree	Rec blocks from last common ancestor	r/m	Genome length (bp)	Bases in clonal frame (bp)	% of recombined genome
<b>1</b>	16	9	8.407	1,627,573	1,588,753	2.99
	7	0	0.978	1,609,340	1,606,957	0.15
<b>2</b>	0	0	0.000	1,666,864	1,666,864	0.00
	7	7	3.923	1,665,656	1,659,185	0.39
<b>3</b>	167	15	35.708	1,617,708	1,220,734	24.90
	162	10	38.985	1,603,282	1,201,817	25.84
<b>4</b>	6	0	10.145	1,616,057	1,593,221	1.41
	15	5	16.141	1,507,717	1,468,305	3.14
<b>5</b>	7	0	5.457	1,651,605	1,614,910	2.22
	52	45	20.475	1,659,211	1,595,168	7.09
<b>6</b>	162	7	46.548	1,621,038	1,225,532	24.71
	161	6	38.660	1,613,935	1,209,402	25.67
<b>7</b>	163	6	44.503	1,613,236	1,215,532	25.05
	171	14	54.284	1,618,186	1,220,026	24.94
<b>8</b>	180	28	39.419	1,621,543	1,223,908	24.84
	160	8	37.994	1,615,211	1,210,318	25.65
<b>9</b>	191	191	6.901	1,579,494	1,287,491	28.63
	172	172	37.839	1,616,244	1,203,711	25.92
<b>10</b>	192	5	143.956	1,616,455	1,209,888	25.48
	197	10	152.767	1,598,501	1,190,972	26.38
<b>11</b>	33	1	48.051	1,648,672	1,514,755	8.12
	32	0	47.809	1,649,085	1,515,390	8.11
<b>12</b>	2	0	4.321	1,653,346	1,647,288	0.37
	12	10	12.071	1,652,455	1,617,482	2.14
<b>13</b>	11	10	9.962	1,582,966	1,569,538	0.85
	36	35	19.288	1,633,910	1,600,539	9.07
<b>14</b>	16	16	31.305	1,645,474	1,588,868	3.44
	32	32	47.809	1,642,338	1,508,643	8.14
<b>15</b>	70	40	63.266	1,655,982	1,494,469	11.56
	35	5	46.657	1,647,567	1,545,999	6.53
<b>16</b>	23	23	7.188	1,651,444	1,521,539	7.87
	23	23	5.289	1,658,746	1,600,962	4.35
<b>17</b>	25	15	50.581	1,617,864	1,541,049	5.10
	72	60	59.602	1,618,755	1,514,712	6.58
<b>18</b>	218	0	42.182	1,597,971	1,119,919	31.12
	306	88	102.167	1,618,051	1,093,701	40.07
<b>SD</b>	85	43	35.448	30,559	189,708	11.74
<b>Mean</b>	87	25	37.240	1,625,375	1,427,987	13.30

SD = standard deviation

Table 7: Recombination information for each isolate in each *C. jejuni* pair considered for rate calibration

Pair no.	Rec blocks from root of subtree	Rec blocks from last common ancestor	r/m	Genome length (bp)	Bases in clonal frame (bp)	% of recombined genome
<b>1</b>	1	0	6.389	1,639,229	1,627,395	0.72
	1	0	6.389	1,634,719	1,622,885	0.72
<b>2</b>	6	6	3.042	1,618,681	1,606,663	0.74
	0	0	0.000	1,618,074	1,618,074	0.00
<b>3</b>	4	0	3.378	1,632,540	1,623,184	0.63
	5	1	0.058	1,626,877	1,618,438	0.66
<b>4</b>	1	0	6.389	1,635,771	1,623,937	0.72
	1	0	6.389	1,597,308	1,585,474	0.74
<b>5</b>	14	14	9.837	1,637,250	1,630,577	2.00
	2	2	2.247	1,639,143	1,637,492	0.10
<b>6</b>	23	23	2.992	1,622,009	1,610,144	1.73
	5	5	0.191	1,613,353	1,602,784	0.66
<b>7</b>	0	0	0.000	1,565,731	1,565,731	0.00
	4	4	0.502	1,605,399	1,576,627	1.79
<b>8</b>	6	6	1.295	1,605,876	1,587,041	1.17
	16	16	1.765	1,602,144	1,595,123	1.62
<b>9</b>	0	0	0.000	1,624,255	1,618,823	1.66
	2	2	0.594	1,635,479	1,630,084	0.33
<b>10</b>	0	0	0.000	1,638,978	1,638,978	0.00
	9	9	2.581	1,617,064	1,597,375	1.34
<b>11</b>	0	0	0.000	1,637,286	1,637,286	0.00
	1	1	0.070	1,634,549	1,633,954	0.04
<b>12</b>	14	5	3.102	1,631,249	1,569,740	3.77
	11	2	1.059	1,635,689	1,595,743	2.45
<b>13</b>	5	5	5.282	1,624,931	1,600,666	1.49
	0	0	0.000	1,635,860	1,635,860	0.00
<b>14</b>	11	3	7.825	1,531,436	1,515,728	1.03
	11	3	8.203	1,479,611	1,464,823	1.01
<b>15</b>	37	12	11.074	1,577,196	1,503,167	5.23
	27	2	4.535	1,578,646	1,527,929	3.21
<b>16</b>	13	5	3.493	1,629,646	1,594,397	3.68
	8	0	2.246	1,585,000	1,559,014	2.55
<b>17</b>	11	11	4.135	1,622,264	1,614,837	0.46
	0	0	0.000	1,608,872	1,608,872	0.00
<b>18</b>	3	3	0.442	1,589,659	1,585,982	0.23
	7	7	6.710	1,626,126	1,625,934	1.05
<b>19</b>	22	11	13.058	1,633,924	1,575,472	3.66
	28	17	10.982	1,571,410	1,495,992	5.30
<b>20</b>	8	8	32.016	1,529,278	1,290,989	18.64
	49	19	35.641	1,600,590	1,336,645	17.51
<b>SD</b>	11	6	7.585	35,024	74,747	3.94
<b>Mean</b>	9	5	5.098	1,609,328	1,579,746	2.22

SD = standard deviation

Table 8: *C. coli* "birthday problem" data and estimates of coalescence across sample time frame

Time cut-off (t)	No. of SNPs in time (st)	No. of pairs < st (X)	Effective lineages (Z)
1	77	8	22,575
2	155	12	15,050
3	232	16	11,288
4	309	21	8,600
5	386	27	6,689
6	464	33	5,473
7	541	41	4,405
8	618	57	3,168
9	696	60	3,010
10	773	64	2,822
11	850	68	2,656
12	928	71	2,544
13	1,005	76	2,376
14	1,082	84	2,150
15	1,159	88	2,052
16	1,237	92	1,963
17	1,314	94	1,921
18	1,391	102	1,771
19	1,469	107	1,688
20	1,546	108	1,672
21	1,623	112	1,613
22	1,700	115	1,570
23	1,778	117	1,544
24	1,855	119	1,518
25	1,932	123	1,468
26	2,010	124	1,456
27	2,087	128	1,411
28	2,164	130	1,389
29	2,241	133	1,358
30	2,319	134	1,348
31	2,396	136	1,328
32	2,473	139	1,299
33	2,551	140	1,290
34	2,628	140	1,290
35	2,705	142	1,272
36	2,783	142	1,272
37	2,860	143	1,263

**mut + rec SNPs/yr rate for *C. coli* = 77.3**  
**Y (all potential pairs) = 180,600**

X = possible pairs; Y = potential pairs; Z = Coalescences ("birthdays")



Table 9: *C. jejuni* "birthday problem" data and estimates of coalescence across sample time frame

Time cut-off ( <i>t</i> )	No. of SNPs in time ( <i>st</i> )	No. of pairs < <i>st</i> ( <i>X</i> )	Effective lineages ( <i>Z</i> )
1	14	76	21,888
2	28	129	12,895
3	42	150	11,090
4	56	163	10,205
5	71	174	9,560
6	85	182	9,140
7	99	190	8,755
8	113	194	8,575
9	127	198	8,401
10	141	209	7,959
11	155	215	7,737
12	169	221	7,527
13	183	225	7,393
14	197	230	7,233
15	212	237	7,019
16	226	242	6,874
17	240	249	6,681
18	254	254	6,549
19	268	256	6,498
20	282	258	6,448
21	296	262	6,349
22	310	267	6,230
23	324	270	6,161
24	338	276	6,027
25	353	278	5,984
26	367	281	5,920
27	381	287	5,796
28	395	292	5,697
29	409	295	5,639
30	423	300	5,545
31	437	303	5,490
32	451	310	5,366
33	465	312	5,332
34	479	314	5,298
35	494	318	5,231
36	508	320	5,198
37	522	324	5,134
38	536	329	5,056
39	550	333	4,995
40	564	335	4,966
41	578	340	4,893
42	592	341	4,878
43	606	342	4,864
44	620	344	4,836
45	635	349	4,766
46	649	352	4,726

**mut + rec SNPs/yr rate for *C. jejuni* = 14.1**  
**Y (all potential pairs) = 1,663,488**

X = possible pairs; Y = potential pairs; Z = Coalescences ("birthdays")