1 **Genomic plasticity and rapid host switching promote the evolution of**

2 **generalism in the zoonotic pathogen *Campylobacter***

3

4 Dan J. Woodcock[1], Peter Krusche[1], Norval J. C. Strachan[2], Ken J. Forbes[3],

5 Frederick M. Cohan[4], Guillaume Méric[5], Samuel K. Sheppard[5,6]∗

6

7 [1]Warwick Systems Biology Centre, Coventry House, University of Warwick, Coventry,

8 CV47AL, UK; [2]School of Biological Sciences, University of Aberdeen, Cruickshank

9 Building. St Machar Drive, Aberdeen, AB24 3UU, UK; [3]School of Medicine and

10 Dentistry, The University of Aberdeen, Foresterhill, Aberdeen, AB25 2ZD, UK;

11 [4]Department of Biology, Wesleyan University, Middletown, CT 06459-0170, USA; [5]The

12 Milner Centre for Evolution, Department of Biology and Biochemistry, University of

13 Bath, Claverton Down, Bath. BA2 7AY, UK; [6]Department of Zoology, University of

14 Oxford, South Parks Road, Oxford, OX1 3PS, UK.

15

16 * Corresponding author: s.k.sheppard@bath.ac.uk;

17

18 Keywords: genome evolution, selection, transmission ecology, adaptation,

19 recombination.

20

21

## Abstract

Horizontal gene transfer accelerates bacterial adaptation to novel environments, allowing selection to act on genes that have evolved in multiple genetic backgrounds. This can lead to ecological specialization. However, little is known about how zoonotic bacteria maintain the ability to colonize multiple hosts whilst competing with specialists in the same niche. Here we develop a stochastic evolutionary model and show how genetic transfer of niche specifying genes and the opportunity for host transition can interact to promote the emergence of host generalist lineages of the zoonotic bacterium *Campylobacter*. Using a modelling approach we show that increasing levels of recombination enhance the efficiency with which selection can fix combinations of beneficial alleles, speeding adaptation. We then show how these predictions change in a multi-host system, with low levels of recombination, consistent with real *r/m* estimates, increasing the standing variation in the population, allowing a more effective response to changes in the selective landscape. Our analysis explains how observed gradients of host specialism and generalism can evolve in a multihost system through the transfer of ecologically important loci among coexisting strains.

## Introduction

40

41

42    Adaptation is typically thought to lead to gradual ecological specialization, in which

43    populations progress towards an optimal phenotype. This can occur among competing

44    organisms in sympatry, particularly when resources diversify or if the cost of maintaining

45    homeostasis in different environmental conditions is high (Van Tienderen 1991).

46    However, resource or host generalism are also widely observed in nature (Fried et al

47    2010, Kassen 2002, Woolhouse et al 2001) and it is generally accepted that natural

48    environmental heterogeneity can promote the maintenance of phenotypic variation where

49    it confers an ecological advantage (Kassen 2002, Van Tienderen 1991).

50

51    Host generalism is exhibited by some bacteria that infect multiple hosts resulting in

52    important implications for the spread of disease from animals to humans. In some

53    zoonotic bacteria, such as *Staphylococcus aureus,* livestock-associated lineages are

54    largely host-restricted (Fitzgerald 2012), allowing the direction and time-scale of host

55    transfer to be estimated by comparison of genotype information (Lowder et al 2009). In

56    contrast, in *Escherichia coli* it is difficult to link ecological niche with genotype as

57    isolates from all major phylogroups are represented in multiple isolate sources (Meric et

58    al 2013). Other organisms including *Campylobacter jejuni* and *Salmonella enterica*

59    (Baumler and Fang 2013) represent an intermediate between these. For example,

60    comparison of *C. jejuni* isolates from various sources, by multilocus sequence typing

61    (MLST) (Sheppard et al 2011) and whole genome sequencing (Dearlove et al 2015,

62    Sheppard et al 2014), has shown evidence for host restricted lineages, found

63    predominantly in one host species, as well as lineages commonly isolated from hosts as

64    diverse as chickens, cattle and wild birds (Gripp et al 2011).

65

66    Out-competition by host specialists, might be expected to select against generalist

67    *Campylobacter* lineages. However, generalists remain among the most common lineages

68    in agricultural animals and are a major cause of human disease (Sheppard et al 2009).

69    There are several factors that may be involved in the emergence of generalism as a

70    successful strategy. First, the development, industrialisation and globalisation of livestock

71    farming have created a vast open niche in which *Campylobacter* has expanded from pre-

72    agriculture wild animal hosts. Second, factors associated with livestock husbandry and

73    habitations promote close contact between different livestock species providing

74    opportunities for *Campylobacter* transmission from one host to another. Third,

75    *Campylobacter* is a highly recombinogenic organism (Wilson et al 2009) and lineages

76    regularly acquire genetic elements through horizontal gene transfer (HGT). This genomic

77    plasticity has the potential to introduce DNA segments, or whole genes, into the recipient

78    genome, potentially conferring novel function.

79

80    The coexistence of *Campylobacter* lineages with generalist and specialist ecology

81    remains poorly understood and little is known about the factors that promote the

82    emergence and maintenance of lineages with distinct ecology (ecotypes). Here, using

83    information on natural genomic variation in 130 *C. jejuni* genomes from chicken and

84    cattle (Sheppard et al 2013b), as well as a computational model of bacterial evolution, we

85    find that the maintenance of genetic variance of putative niche-specifying genes in the

86    population can promote the emergence of generalism as seen in nature. By quantifying

87    how resource competition, rapid host switching, and horizontal gene transfer interact to

88    affect the variance in the population, we provide a generalised framework for considering

89    the emergence of generalist ecotypes.

90

91    **Material and methods**

92

93    **Bacterial genomes**

94    A total of 130 *C. jejuni* and *C. coli* isolates were used in this study, including 87

95    representative strains sampled from chicken and cattle (Table S1). The genomes were

96    previously published and isolates were described as belonging to clonal complexes on the

97    basis of sharing 4 or more identical alleles at seven multilocus sequence typing loci

98    (Sheppard et al 2013a, Sheppard et al 2013b). In total for this study, 30 genomes from the

99    ST-21 clonal complex, 28 from the ST-45 clonal complex, 7 from the ST-353 clonal

100   complex, 6 from the ST-206 clonal complex, 6 from the ST-61 clonal complex, and 5

101   from the ST-48 clonal complex were used. ST-353 complex isolates are known to be

4

102  chicken-associated and ST-61 complex isolates to be cattle-associated, whereas ST-21,

103  ST-45, ST-206 and ST-48 are host generalists (Sheppard et al 2014). Genomes were

104  archived on a web-based platform system based on BIGSdb, which also implements

105  analytic and sequence exporting tools (Jolley and Maiden 2010). An additional 75

106  genomes representing 74 different STs from *C. jejuni* and *C. coli* were used. These

107  genomes were sequenced as part of other studies (Food Standards Scotland i-CaMPS-3

108  Contract S14054, DEFRA project OZ0625), and from the PubMLST Database (Jolley

109  and Maiden 2010).

110

111  **Model input data**

112  Using a gene-by-gene approach (Maiden et al 2013, Sheppard et al 2012), loci in the 130

113  genomes were identified by BLAST comparison to the *C. jejuni* strain NCTC11168

114  reference genome (Genbank accession number: AL111168) with a >70% nucleotide

115  sequence identity on ≥50% of sequence considered sufficient to call a locus match, as in

116  other studies (Meric et al 2014, Meric et al 2015, Pascoe et al 2015). A whole-genome

117  MLST (Maiden et al 2013, Sheppard et al 2012) matrix was produced summarizing the

118  presence/absence and allelic diversity at each locus in each genome, based upon these

119  BLAST parameters. From this matrix, 1080 core genes were found to be shared by all

120  cattle and chicken isolates from our dataset. The proportion of each allele at each locus

121  was calculated in both cattle and chicken and then subtracted to identify alleles that were

122  common in one group, but rare in the other. These were then summed these values at

123  each locus to get the discriminative capacity of each locus for each host species. Loci in

124  which the alleles segregated by host were considered as a proxy for niche-specifying

125  genes in the model. While this was done in preference to simulating data, no inference is

126  made based on the function or potential adaptive advantage conferred. A total of 5

127  putative niche-specifying genes per host (chicken and cattle) and 5 MLST genes picked

128  at random (15 genes in total) were used as the input genotype for the model (**Table S2**).

129  The inclusion of MLST loci (*aspA*, *uncA*, *pgm*, *glnA*, *gltA*) provides a reference for

130  comparison.

131

132  **The Genome Evolution by Recombination and Mutation (GERM) model**

133      The GERM model is a simplified representation of a bacterial population and associated

134      processes, which allows us to simulate bacterial evolution *in silico* by tracking individual

135      bacteria of variable genotypes as they are exposed to various selective environmental

136      pressures. Furthermore, by simulating with a stochastic sampling algorithm, we can also

137      incorporate some degree of the randomness inherent in natural populations, and hence

138      investigate the importance of stochastic effects by performing simulations with the same

139      initial population. Similar models have been proposed previously (Levin and Cornejo

140      2009), and our approach extends and builds on these models in a number of ways. In the

141      model in this study, each individual bacterial cell is represented as a 15-locus genotype as

142      described above. In a population of $N$ cells, each cell $C$ is described entirely by the

143      alleles $a$ at locus $j$ which compose the genome, denoted $a_j$, $j = 1, 2, ..., 15$. Each allele at

144      each locus is represented by an integer ranging from 0 to $\infty$. Basing our algorithm on real

145      data, where there are approximately 20 possible alleles at each locus ($20^{15} \approx 3.27 \times 10^{19}$

146      different combinations) presents computational challenges if we model the population

147      using proportions of genotypes as in previous models (Levin and Cornejo 2009). To

148      account for this we store the entire population at any one time and perform operations at

149      the individual level. Working with the population directly, instead of adjusting

150      proportions of STs, allows the investigation of the population dynamics at different

151      population sizes, which is particularly pertinent after selective sweeps when the number

152      in the population itself will drop as the population adapts to the new environment. The

153      model incorporates six basic processes: mutation, recombination, resource consumption,

154      cell death, cell division and host migration. Each process occurs once per generation and

155      is stochastic, therefore occurring with a probability defined for each cell. These can be

156      interpreted as rates per generation.

157

158      **Cell division, mutation and recombination**

159      Mutation and recombination occur at the level of the individual locus and cell death and

160      cell division occur at the individual bacterial cell level. With cell division, an identical

161      copy of the cell that divides is added to the population, this occurs with probability $b$.

162      Mutation occurs with probability $m$ and, unlike existing models (Levin and Cornejo

163      2009) any allele that mutates is deemed to offer no selective advantage (the fitness of that

164  allele is 0). Similarly, a recombination event occurs at probability $r$ and an allele that

165  recombines is assigned the value of another allele randomly chosen from those at the

166  same locus within the current population. In natural systems, these processes are typically

167  considered rare with upper rate estimates for homologous recombination of $10^{-6}$ per gene

168  per generation (Wiedenbeck and Cohan 2011). However, these rate estimates can be

169  affected by a number of factors (Barrick and Lenski 2013, Vos and Didelot 2009) and so

170  typically in studies of bacterial evolution, the preferred measure of the magnitude of

171  recombination is the relative frequency of recombination compared to mutation, the $r\,/\,m$

172  ratio (Falush et al 2001, Fearnhead et al 2005, Feil 2004, Fraser et al 2005, Milkman and

173  Bridges 1990). To quantify the effects of varying levels of homologous recombination on

174  niche adaptation in the GERM model, we used $r/m$ ratio as the ratio of rates at which

175  alleles are substituted as a result of recombination and mutation. Partly for reasons of

176  computational tractability, the GERM model simulates a simplification of the size and

177  complexity of a natural system, and imposes enhanced selection against maladapted

178  sequences types. Because of this, $r$ and $m$ rate estimates from natural populations are

179  adjusted so that sequence types that arise in the population have a similar opportunity on

180  average to proliferate in both the model and the natural environment and are not removed

181  from the population by chance alone. Consistent with existing estimates from multiple

182  bacterial species including *Campylobacter* (Vos and Didelot 2009)*,* we run simulations at

183  $r\,/\,m$ ratios ranging from 0 to 100 corresponding to a mutation rate of 0.01, with

184  recombination rates of 0, 0.001, 0.01, 0.1 and 1 to facilitate comparisons with natural

185  populations.

186

187  **Fitness**

188  In this study, host-specific alleles at niche specifying genes are considered to confer a

189  fitness advantage to the cell in one or other host. The fitness of allele $a$, in a given host

190  $h$ is defined as $f^{\{h\}}(a)$ and this reflects the fitness conferred by that allele to its

191  environment, with 1 corresponding to a perfectly adapted allele conferring maximal

192  fitness, ranging to an allele that provides no benefit to the survival of the cell and has 0

193  fitness. We then follow Levin (Levin and Cornejo 2009) and calculate the fitness of an

194  individual cell as the sum of the allelic fitness values assuming that each allele

195     contributes equally. The fitness function is therefore: $F(C) = \dfrac{1}{n} \sum_{j=1}^{15} f^{\{h\}}(a_j)$. As a

196     consequence, a different fitness landscape, defined by another host, affects the fitness of a

197     cell through the fitnesses of its constituent alleles. It is by changing the fitness

198     landscapes that host transition is simulated. Linkage disequilibrium is not factored into

199     the model, although genes were selected from divergent genomic positions to limit

200     linkage effects, and the model was not designed to test the specific function of individual

201     genes. In nature, complex interactions between genes and the environment are likely to

202     correspond to fitness function involving operations somewhere between additive and

203     multiplicative (Phillips 2008) but less information on this is available for bacteria than in

204     the more well studied diploid scenario.

205

206     **Resource**

207     The fate of a bacterium in the GERM model is not only dependent on fitness, but also on

208     the availability of resources. This introduces a dynamic relationship between the fate of

209     cells, their fitness, and the population size, and confers a soft carrying capacity for a

210     niche dependence on the interplay between these aspects. Resource is modeled as a

211     generic entity for which no distinction is made for the type of resource in a given niche,

212     and the affinity for the consumption of a resource by a bacterium is independent of

213     genotype. For each cell, the chance of using a resource occurs with probability $u^+$,

214     multiplied by the amount of available resource. If a cell is already using a resource, it

215     finishes with probability $u^-$ in which case the resource is then consumed and is not

216     returned to the environment. A resource is generated with constant probability $g$ and is

217     added to the pool of available resource. As such, for any given death in the population,

218     there is a probability that it is caused either by a lack of fitness, or the inability to find and

219     utilize resource. For a given fitness, the chance of death due to fitness or resources

220     changes as the number of cells rises. For very small populations, the cause of death is

221     predominantly fitness-related as resources are abundant. As the population increases, and

222     the chance of finding a resource diminishes, resource-related death quickly starts to

223     influence the fate of the cell but then plateaus as the population approaches carrying

224     capacity. This nonlinear relationship is because, when unconstrained, the population will

225     double in size during the course of a generation, where the resources will only increase

226     by a fixed amount each time, regardless of population size. Conversely, as fitness

227     increases, the proportion of deaths caused by fitness decreases to the point that resource

228     related death becomes dominant. This is important as the mechanism of death is different

229     in each case: in fitness-related death, the probability of death is inversely proportional to

230     the fitness. However, for resource-related death, the chance of death is independent of

231     fitness and occurs entirely at random. As such, when resources become scarce, the

232     fitness of the members of a population becomes largely less important to their fate.

233     Therefore, resource availability is fundamental to the population dynamics, as well as

234     suppressing the speed of convergence to the optimal genotype and maintaining genetic

235     variance.

236

237     **Cell death**

238     Cell death occurs in two stages, fitness-related death and resource-related death. Fitness-

239     related death is dependent on a fitness function, as discussed above and detailed further

240     the supplementary material, which provides a probability of dying per generation for a

241     given genotype. Resource-related death can happen only when a cell is not using a

242     resource and occurs with probability $d$ regardless of fitness. In both cases, the cell is

243     removed from the population. The algorithm also incorporates a host transition or a

244     selective sweep by switching between fitness functions. In this case, the resources can

245     also be reset (although those already using a resource will continue as before). Further

246     details of the model and stochastic simulation are included in the supplementary material.

247

248     **Simulation**

249     To account for the randomness inherent in the process of evolution, we simulated the

250     system with a variant of a stochastic sampling algorithm (SSA) (Gillespie 1977),

251     involving Monte Carlo decision processes, was used for simulating the system. The

252     number of cells required to represent a natural population is extremely large, therefore

253     some constraints were imposed to ensure computational tractability. We sampled from

254     Poisson distributions and binomial distributions when the quantities in question were

255     unbounded and bounded respectively. For simplicity, we set the probability of a cell

256    division event $b$ to be equal to 1 so that every iteration can be interpreted as a generation.

257    For mutation and recombination we imposed a constraint that each of these events can

258    occur only once every time step. There were ten time steps per generation. The algorithm

259    iterated through the five processes in order of resource allocation, mutation,

260    recombination, cell death and cell division. We use three fitness functions, one each

261    from chicken and cow, which we subsequently refer to as "Host 1" and "Host 2"

262    respectively, and one derived using data from both hosts referred to as the composite

263    fitness function. The composite fitness function is used to initialize the simulation and

264    allow a burn-in period where the genetic composition to arrange themselves in a way that

265    favors neither host. The algorithms start in using the composite fitness function to allow

266    the algorithm to burn-in and deplete unfit allele combinations without purging alleles

267    from one of the hosts. When the population has recovered and reached equilibrium with

268    resource generation, we switch the fitness function to one of the single hosts. This

269    occurred after 200 generations. We subsequently alternated the fitness function between

270    chicken and cattle with the frequency defined by the user. This corresponds to the

271    movement of the entire population to the new host, a simplification of the more realistic

272    process in which bacteria would move to one of several distinct and concurrently

273    evolving ecosystems. This layer of abstraction was partially chosen for computational

274    tractability, but also to aid in the interpretation of results as the number of stochastic

275    events would increase exponentially as the model complexity increases which would

276    potentially dominate any patterns in our results. As such, this constraint is suitable for

277    the scope of this study, particularly as our model is predicated on *Campylobacter*

278    populations which are often transmitted *en masse* through stool. To initialize the

279    algorithm, we generated a population of 50 million cells consisting of alleles found in the

280    data set, randomly assigned throughout the population in a uniform manner rather than

281    weighted by their abundance in the data set. The same initial population was used for

282    every simulation run. This way, the populations and processes can be kept identical, and

283    so that the only differences between runs with the same parameters were stochastic

284    effects.

285

286    Two experiments were carried out at a range of recombination to mutation ratios (0-100).

287    1.  Long term adaptation followed by host switch – after the first composite burn in,
288        a switch to host 2 was performed. This was run for 1000 generations to simulate
289        long term colonization and adaptation and was then followed by a switch to host
290        1.
291    2.  Rapid host switching – after the composite burn in – rapid host switching (every
292        200 generations) was performed.
293

294    **Host generalism and phylogeny**
295    We used 75 genomes from isolates representing 74 *C. jejuni* and *C. coli* STs, from 30239
296    pubMLST isolate records (http://pubmlst.org/campylobacter/), to investigate the degree
297    of host generalism among different lineages. From these, concatenated gene-by-gene
298    alignments of 595 core genes (Sheppard et al 2013b), constructed as previously published
299    (Meric et al 2014, Meric et al 2015), were used to infer phylogenetic relationships. A tree
300    was constructed using the neighbor-joining algorithm and each ST was labelled with the
301    number of distinct hosts from which it has been isolated, based on data submitted to the
302    pubMLST database. The 74 STs were reported to have been isolated from 20 distinct
303    animal hosts, including humans. The maximum number of hosts associated with a single
304    ST was 12 (ST-45) and the minimum was 1 (ST-58). A tree and the heatmap representing
305    the number of hosts, was prepared and visualized in Evolview (Zhang et al 2012). The
306    mean host species richness score was correlated with the genetic distance (derived from
307    the number of SNPs) from the tip of the tree to the first branching point for each of the
308    isolates.
309

310    **Results**
311

312    **Long term adaptation promotes specialism but recombination enhances colonization**
313    **in a subsequent host transition**
314    The effect of homologous recombination was characterized in 200 independent
315    population GERM simulations, with a transition from the composite niche to Host 1 after
316    200 generations, followed by a transition to Host 2 after another 1000 generations.
317    Simulations were performed at five $r/m$ ratios: 0; 0.1; 1; 10; 100. In all simulations, the

318    mean number of cells decreased sharply from the initial condition and then recovered to

319    approach an equilibrium level between birth and death just before the transition to Host 1

320    (Figure 1). In the composite niche, level of proliferation was proportional to the rate of

321    recombination with concomitant increase in mean fitness and population genetic variance

322    (Figures 1C and 1D). This is because higher recombination rates result in greater genetic

323    variance, and so by Fisher's fundamental theorem of natural selection (Fischer 1930), the

324    rate of increase in the mean fitness will be greater and the population will thrive.

325    Following the first host transition (Figure 1A, numeral I), there was a brief increase in the

326    population due to increased resource availability, but in all cases the mean number of

327    cells quickly returned to the equilibrium state where it continued until the next host

328    transition (numeral II), with the mean number of cells ordered largely as before, with the

329    number of cells increasing with recombination rate (Figure 1C). After the transition to

330    Host 2, the populations were decimated in all cases, as the alleles which would have

331    conveyed enhanced fitness in Host 2 have been purged from the population, meaning that

332    the bacteria are ill-equipped to survive in the new host and so they will die. Almost all of

333    the populations died out in the low recombining groups ($r/m$=0 and 0.1), a large

334    proportion died out in the intermediate group ($r/m$=1), with the greatest number of

335    surviving populations in the highly recombining groups.

336

337    **Intermediate recombination rates enhance population mean fitness after multiple**

338    **rapid host transitions**

339    As in the single host transition model, the mean number of cells in the composite niche

340    reached equilibrium levels ordered by their recombination rate (Figure 2A) and consistent

341    with their mean fitness levels (Figure 2B). At the end of growth in the composite niche,

342    the intermediate and high recombination rates ($r/m$=1, 10 and 100) have a similar ability

343    to survive, with the highest recombination level displaying high variance (**Figure 2C**).

344    After the composite niche a number of host transitions were simulated where the mean

345    number of cells shifted between two equilibrium states depending on the host species.

346    The mean number of cells was always higher in Host 1 because some alleles conferring

347    increased fitness in Host 2 will inevitably be purged from the populations in the Host 1

348    niche. Reversing the species order had the same effect for Host 2.  At the end of the last

349   Host 2 growth cycle all the recombining populations ($r/m$>0) had a similar mean number

350   of cells, but the variance differed, with the smallest variance at $r/m$=1 and $r/m$=10.

351   Similarly, in the final Host 1 niche, we can see in Figure 2A that the intermediate

352   recombination rate ($r/m$=1) is associated with the highest mean number of cells. This

353   shows that in contrast to the single host transition model, recombination at a high level is

354   not advantageous under repeated host switching.

355

356   **Emergence of host generalist strategy as a consequence of frequent host switches**

357   We used a Dirichlet process clustering algorithm (Kurihara et al 2007) on all simulations

358   to identify characteristic profiles of population dynamics for the different recombination

359   rates (Figure 3). Three broad population dynamic profiles were observed: (i) populations

360   that were primarily adapted to Host 1 (Clusters 1-5); (ii) populations that were primarily

361   adapted to Host 2 (Clusters 7-9); (iii) populations that were adapted to both niches

362   (Cluster 6). This is consistent with a classification as a specialist for either species, or as a

363   generalist.  The membership of these 3 adaptation profile types relates to recombination

364   rate (Table 1). Simulations with no recombination were predominantly found in Cluster

365   2, with a substantial amount found in other clusters. This is to be expected as the outcome

366   was driven entirely by stochasticity acting on the population and so genes were purged

367   almost at random in the composite niche, yielding a set of outcomes which were

368   maintained during the host transitions as more alleles were lost. In contrast, it can be seen

369   that simulations of all of the recombining populations are predominantly found in Cluster

370   4, which is a Host 1 specialist cluster, albeit with a relatively high equilibrium population

371   number during the Host 2 niche compared to the other Host 1 specialist clusters. In

372   Cluster 4, it can be seen that an $r/m$=1 gives the greatest occupancy, at 91.8%, explaining

373   the high numbers of cells seen in Figure 2. The membership of the generalist cluster,

374   Cluster 6, is also represented across all recombination rates, with the highest percentage

375   coming from a relatively low recombination rate ($r/m$=0.1).

376

377   **The gradient of host generalism is mirrored in natural *Campylobacter* populations**

378   The degree of host specialism and generalism in *in silco*, resulting from model

379   simulations, was compared to data from natural *Campylobacter* populations. Mapping

13

380   isolation source data of 30239 *Campylobacter* STs within the PubMLST database (Jolley

381   and Maiden 2010) to a core genome neighbor joining tree of the 74 common *C. jejuni*

382   and *C. coli* STs revealed that few STs demonstrate absolute generalism or specialism

383   (Figure 4a). Rather there is a gradient ranging from STs predominantly found in one host

384   to those frequently isolated from multiple host sources (Figure 4b). The clustering of

385   isolate pairs, with shared host source richness, was estimated by correlating nucleotide

386   divergence in the core genome with the number of hosts (Figure 4c) and by a

387   randomisation/ permutation test which showed p<0.000001 (data not shown). STs on the

388   phylogeny were located close to STs with similar host species richness suggesting that

389   there is some evolutionary signal which determines the likelihood of an ST's degree of

390   specialism.

391

392   **Discussion**

393

394   The GERM model provides a context for considering how genome plasticity may

395   influence the proliferation of *Campylobacter* in a multihost environment. In simple

396   simulations, rapid acquisition of niche-specifying genes promoted better colonization in a

397   new host. This is consistent with the Fisher–Muller evolutionary model (Fisher 1930,

398   Muller 1932) where recombination functions to bring together fit alleles, which would

399   otherwise compete for fixation in the population, into a single lineage speeding the

400   overall increase in population mean fitness (Gerrish and Lenski 1998). In line with

401   classical population genetic theory for sex (Barton and Charlesworth 1998, Felsenstein

402   1974, Weismann 1904), the efficiency of selection on bacteria is enhanced by this

403   shuffling of alleles. Therefore, the population with the highest recombination rate will

404   expand to fill the niche more rapidly after a genetic bottleneck. This demonstrates a clear

405   short-term adaptive benefit to rapid recombination, but does not explain why most

406   bacteria recombine at low rates in nature.

407

408   Where survival is predominantly influenced by a few genetic determinants, for example

409   the acquisition of essential antibiotic resistance genes (Spratt et al 1989), high

410   recombination rates would be favoured. However, this is an unusually simple

411  evolutionary scenario and bacterial habitats comprise numerous interacting selective
412  pressures. Because increased genetic variation leads to faster adaptation (Fisher 1930),
413  the potential for the population to survive future genetic bottlenecks is related to the
414  fitness variance. In populations that recombine at a low rate, a relatively high fitness
415  variance is often maintained. Therefore, if the species is likely to encounter frequent
416  environmental changes, such as host switches, it may be beneficial to have a lower
417  recombination rate than would be optimal in the Fisher-Muller model.

418

419  In nature, each niche will have an associated carrying capacity. As a species expands to
420  approach this carrying capacity, competition will inevitably increase, meaning that the
421  influence of external factors will be of increased importance to the fate of the organism.
422  To avoid this competition, an organism could adapt to a new less competitive niche,
423  consistent with the 'tangled bank hypothesis' (Bell 1982, Doebeli 1996, Koella 1988) for
424  the evolution of sexual reproduction. In our competitive model, fitness variance is higher
425  in populations with low recombination rates, providing a competitive advantage under the
426  tangled bank hypothesis that would facilitate a transition to a new niche. This provides an
427  explanation for the low recombination rates observed in some natural bacterial
428  populations, and explains the presence of multiple lineages as infrequent recombination
429  will allow the uptake of adaptive genes but may be too infrequent to prevent adaptive
430  divergence between lineages (Wiedenbeck and Cohan 2011).

431

432  Based on model simulations, the most favorable recombination to mutation ratio for
433  promoting *Campylobacter* survival in the new niche whilst maintaining fitness variance
434  within the population was $r/m$ =0.1-1 which is comparable to that calculated in natural
435  *Campylobacter* populations ($r/m$ =0.44) (13). In multiple niche transition simulations,
436  at intermediate recombination rates ($r/m = 1$), many populations did not completely die
437  out, but resisted the introduction to a novel host recovering after a few passages, as seen
438  by population size and mean fitness increases over time. This provides a context for
439  considering the balance between rapid adaptation, mediated by recombination, and
440  maintenance of genetic variance, allowing each population to survive in both host
441  environments over time (Figure 1 and Figure 2). Absolute host specialists, which went

442   extinct in the second host, were uncommon, and most populations demonstrated some

443   capacity to survive in both hosts.

444

445   In most cases, simulated genotypes were more successful in one or other niche. This is

446   consistent with evidence from natural populations where lineages such as the ST-257 and

447   ST-61 clonal complexes are predominantly associated with chicken and cattle

448   respectively but are also isolated from both niches (Sheppard et al 2014). However, in

449   some multihost simulations, populations emerged that were affected very little by the

450   host switches (Figure 3). These populations can be considered true ecological generalists,

451   comparable to the ST-21 and ST-45 clonal complexes that are regularly isolated from

452   chickens, cattle and other hosts (Sheppard et al 2014).

453

454   Ecological specialism and generalism have been well described in animals. For example,

455   the Giant Panda, *Ailuropoda melanoleuca*, is a paradigm of specialism, confined to six

456   isolated mountain ranges in south-central China, where bamboo comprises 99% of its diet

457   (Lü et al 2008), while the American Black Bear, *Ursus americanus*, is a generalist,

458   opportunist omnivore with a broad range including temperate and boreal forests in

459   northern Canada and subtropical areas of Mexico (Garshelis et al 2008). Specialization is

460   a potentially precarious strategy as change to the environment can cause extinction if

461   organisms are unable to move between niches or hosts. Consistent with this, generalist

462   lineages would be expected to be older, preceding specialists which cluster closer to the

463   tips of the phylogenetic tree (Stireman 2005). In *Campylobacter,* there was no correlation

464   between tree tip length and number of hosts suggesting similar evolutionary timescales

465   for specialist and generalist STs. This contrast with studies of metazoans may, in part, be

466   explained by the ability of *Campylobacter* to rapidly acquire niche-specifying elements

467   leading to rapid adaptation of multiple lineages.

468

469   In addition to genomic plasticity, the scale of environmental variation can act on the type

470   of ecological strategy observed among the inhabitants (Futuyma and Moreno 1988,

471   Kassen 2002). For example, in a highly stratified environment, adaptation may occur

472   early with traits becoming fixed, whereas in a graduated environment individuals may be

473  more likely to show a reversible phenotype response (Kassen 2002). Specific niches of

474  different bacterial lineages may be less well defined than for some animal species, but

475  isolation source data from the PubMLST database indicates that *Campylobacter* STs

476  show a gradient of host generalism. This is influenced by the opportunity to colonize the

477  new niche and the capacity to survive the transition. Consistent with the findings in this

478  study, rapid host transitions provide an opportunity for the proliferation of organisms

479  with varying levels of host specialization including generalists that are capable of

480  proliferating in more than one host.

481

482  In this study, by combining data from natural bacterial populations and a selection driven

483  computer modeling, we simulated and predicted the evolution of host association

484  strategies in the zoonotic pathogen *Campylobacter*. In practice, bacteria in natural

485  populations usually exist in complex changeable ecosystems with numerous selection

486  pressures. Here we show that recombination allows a more rapid response after a genetic

487  bottleneck, as in a host transition, by increasing the efficiency with which selection can

488  fix combinations of beneficial alleles. Furthermore, in a dynamic setting of host

489  switching, recombination rate was observed to be a key factor in the colonization and

490  maintenance in multiple niches. Livestock in modern intensive agricultural systems are

491  different to ancestral host populations in numerous ways associated with diet and

492  stocking density. The implications of this are potentially significant as, under conditions

493  favoring rapid host switching, the emergence of host generalist zoonotic pathogens can

494  be simulated. Our model therefore provides a context for considering how recombining

495  bacteria, such as *Campylobacter,* could evolve to meet the challenges of anthropogenic

496  environmental change. This could promote the emergence of multi-host pathogens and

497  increase their capacity to overcome deliberate human interventions.

## Acknowledgements

## Conflict of interest statement

Authors declare no conflict of interest.

## Figure and table legends

**Figure 1. Long term host adaptation followed by a host switch with different recombination rates**. The number of cells (A), population mean fitness (C) and population variance (D) were monitored in 200 independent simulations performed at recombination to mutation ratios of: $r/m$=0 (blue), $r/m$=0.1 (red), $r/m$=1 (green), $r/m$=10 (magenta) and $r/m$=100 (brown). Model parameters were different in phases I, II and III. Phase I corresponds to a composite niche no fitness-related selection. The transition from I to II corresponds to the addition of selective pressure favoring genes specifying adaptation to host 1, and the transition from phase II to III, corresponds to a single host transition, with a change to selective pressures favoring genes specifying adaptation to host 2. Panel B shows the number of cells at the end of every phase for populations with different recombination rates.

**Figure 2. Rapid multiple host transitions with different recombination rates**. Number of cells (A) and mean fitness (B) for recombination to mutation ratios of $r/m$=0 (blue), $r/m$ =0.1 (red), $r/m$ =1 (green), $r/m$ =10 (magenta) and $r/m$ =100 (brown) as they progress through several niche transfers (broken lines), for which selective pressures are alternatively imposed to favor genes specifying adaptation to one or the other host. (C) Population growth distribution at various recombination rates at generations 200

529 (transition from composite niche to the first host), 1400 (last transition growth phase in

530 host 2) and 1600 (end of simulation after 6 host switches).

531

532 **Figure 3. Population dynamic profile clusters**. Population profiles for all simulates at

533 tested recombination rates. Nine distinct profile clusters were identified and the mean

534 number of cells (A) and the mean fitness (B) is shown for example simulations for profile

535 cluster. Black broken lines indicate a host transition where selective pressures switch in

536 favour of the other host. Error bars are given at the midpoint of each niche and

537 correspond to one standard deviation. Ecological groups were inferred from the profiles,

538 with specialist groups for Host 1 (red line) and Host 2 (blue line) as well as a generalist

539 group (black line).

540

541 **Figure 4. Phylogeny of generalist and specialist *Campylobacter* lineages**. (A)

542 Phylogenetic tree and isolation source of 74 common C. jejuni and C. coli sequence types

543 (STs). The tree was constructed from a concatenated gene-by-gene alignment of 595 core

544 genes, using the neighbor joining (NJ) algorithm. The heatmap represents the number of

545 hosts in which particular STs were isolated, based on analysis of 30239 pubMLST isolate

546 records. The scale bar represents the number of substitutions per site; (B) Quantification

547 of the variation in number of hosts for all lineages shown in the tree, highlighting a

548 gradient between host specialism and generalism; (C) Comparison of clustering on the

549 tree, calculated as the estimated number of SNP corresponding to the average branch tip

550 distance to the last common ancestor (LCA) with the number of hosts of isolation of each

551 examined ST. There is no correlation between the two datasets (linear regression;

552 $r^2=0.037$).

553

554 **References**

555

556 Barrick JE, Lenski RE (2013). Genome dynamics during experimental evolution.
557 *Nat Rev Genet* **14:** 827-839.

558

559 Barton NH, Charlesworth B (1998). Why Sex and Recombination? *Science* **281:**
560 1986-1990.

561

562    Baumler A, Fang FC (2013). Host specificity of bacterial pathogens. *Cold Spring*
563    *Harbor perspectives in medicine* **3:** a010041.
564
565    Bell G (1982). *The Masterpiece of Nature: the Evolution and Genetics of*
566    *Sexuality.* Croom Helm: London.
567
568    Dearlove BL, Cody AJ, Pascoe B, Meric G, Wilson DJ, Sheppard SK (2015).
569    Rapid host switching in generalist Campylobacter strains erodes the signal for
570    tracing human infections. *ISME J.*
571
572    Doebeli M (1996). Quantitative Genetics and Papulation Dynamics. *Evolution* **50:**
573    532-546.
574
575    Falush D, Kraft C, Taylor NS, Correa P, Fox JG, Achtman M *et al* (2001).
576    Recombination and mutation during long-term gastric colonization by
577    Helicobacter pylori: estimates of clock rates, recombination size, and minimal
578    age. *Proc Natl Acad Sci U S A* **98:** 15056-15061.
579
580    Fearnhead P, Smith NG, Barrigas M, Fox A, French N (2005). Analysis of
581    recombination in Campylobacter jejuni from MLST population data. *J Mol Evol*
582    **61:** 333-340.
583
584    Feil EJ (2004). Small change: keeping pace with microevolution. *Nat Rev*
585    *Microbiol* **2:** 483-495.
586
587    Felsenstein J (1974). The evolutionary advantage of recombination. *Genetics* **78:**
588    737-756.
589
590    Fischer R (1930). *The Genetical Theory of Natural Selection.* The Clarendon
591    Press: Oxford.
592
593    Fisher RA (1930). *The genetical theory of natural selection.*, 1 edn, vol. 1.
594    Clarendon Press: Oxford.
595
596    Fitzgerald JR (2012). Livestock-associated Staphylococcus aureus: origin,
597    evolution and public health threat. *Trends Microbiol* **20:** 192-198.
598
599    Fraser C, Hanage WP, Spratt BG (2005). Neutral microepidemic evolution of
600    bacterial pathogens. *Proc Natl Acad Sci U S A* **102:** 1968-1973.
601
602    Fried G, Petit S, Reboud X (2010). A specialist-generalist classification of the
603    arable flora and its response to changes in agricultural practices. *BMC ecology*
604    **10:** 20.
605
606    Futuyma DJ, Moreno G (1988). The Evolution of Ecological Specialization.
607    *Annual Review of Ecology and Systematics* **19:** 207-233.

608

609    Garshelis DL, Crider D, van Manen F, Group ISBS (2008). Ursus americanus.
610    The IUCN Red List of Threatened Species.

611

612    Gerrish PJ, Lenski RE (1998). The fate of competing beneficial mutations in an
613    asexual population. *Genetica* **102-103:** 127-144.

614

615    Gillespie DT (1977). Exact Stochastic Simulation of Coupled Chemical
616    Reactions. *The Journal of Physical Chemistry* **81:** 2340–2361.

617

618    Gripp E, Hlahla D, Didelot X, Kops F, Maurischat S, Tedin K *et al* (2011). Closely
619    related Campylobacter jejuni strains from different sources reveal a generalist
620    rather than a specialist lifestyle. *Bmc Genomics* **12:** 584.

621

622    Jolley KA, Maiden MC (2010). BIGSdb: Scalable analysis of bacterial genome
623    variation at the population level. *BMC Bioinformatics* **11:** 595.

624

625    Kassen R (2002). The experimental evolution of specialists, generalists, and the
626    maintenance of diversity. *Journal of Evolutionary Biology* **15:** 173-190.

627

628    Koella JC (1988). The Tangled Bank: The maintenance of sexual reproduction
629    through competitive interactions. *Journal of Evolutionary Biology* **1:** 95-116.

630

631    Kurihara K, Welling M, Vlassis N (2007). Accelerated variational Dirichlet process
632    mixtures. *Advances in Neural Information Processing Systems* **19:** 761-768.

633

634    Levin BR, Cornejo OE (2009). The population and evolutionary dynamics of
635    homologous gene recombination in bacterial populations. *Plos Genet* **5:**
636    e1000601.

637

638    Lowder BV, Guinane CM, Ben Zakour NL, Weinert LA, Conway-Morris A,
639    Cartwright RA *et al* (2009). Recent human-to-poultry host jump, adaptation, and
640    pandemic spread of Staphylococcus aureus. *Proc Natl Acad Sci U S A* **106:**
641    19545-19550.

642

643    Lü Z, Wang D, Garshelis DL, Group ISBS (2008). Ailuropoda melanoleuca. The
644    IUCN Red List of Threatened Species.

645

646    Maiden MC, van Rensburg MJ, Bray JE, Earle SG, Ford SA, Jolley KA *et al*
647    (2013). MLST revisited: the gene-by-gene approach to bacterial genomics. *Nat*
648    *Rev Microbiol* **11:** 728-736.

649

650    Meric G, Kemsley EK, Falush D, Saggers EJ, Lucchini S (2013). Phylogenetic
651    distribution of traits associated with plant colonization in Escherichia coli. *Environ*
652    *Microbiol* **15:** 487-501.

653

654 Meric G, Yahara K, Mageiros L, Pascoe B, Maiden MC, Jolley KA *et al* (2014). A
655 reference pan-genome approach to comparative bacterial genomics:
656 identification of novel epidemiological markers in pathogenic Campylobacter.
657 *PLoS One* **9:** e92798.

659 Meric G, Miragaia M, de Been M, Yahara K, Pascoe B, Mageiros L *et al* (2015).
660 Ecological Overlap and Horizontal Gene Transfer in Staphylococcus aureus and
661 Staphylococcus epidermidis. *Genome Biol Evol* **7:** 1313-1328.

663 Milkman R, Bridges MM (1990). Molecular evolution of the Escherichia coli
664 chromosome. III. Clonal frames. *Genetics* **126:** 505-517.

666 Muller HJ (1932). Some genetic aspects of sex. *The American Naturalist* **66:**
667 118-138.

669 Pascoe B, Meric G, Murray S, Yahara K, Mageiros L, Bowen R *et al* (2015).
670 Enhanced biofilm formation and multi-host transmission evolve from divergent
671 genetic backgrounds in Campylobacter jejuni. *Environ Microbiol.*

673 Phillips PC (2008). Epistasis--the essential role of gene interactions in the
674 structure and evolution of genetic systems. *Nat Rev Genet* **9:** 855-867.

676 Sheppard SK, Dallas JF, Strachan NJ, MacRae M, McCarthy ND, Wilson DJ *et al*
677 (2009). Campylobacter genotyping to determine the source of human infection.
678 *Clin Infect Dis* **48:** 1072-1078.

680 Sheppard SK, Colles FM, McCarthy ND, Strachan NJ, Ogden ID, Forbes KJ *et al*
681 (2011). Niche segregation and genetic structure of Campylobacter jejuni
682 populations from wild and agricultural host species. *Mol Ecol* **20:** 3484-3490.

684 Sheppard SK, Jolley KA, Maiden MCJ (2012). A Gene-By-Gene Approach to
685 Bacterial Population Genomics: Whole Genome MLST of Campylobacter. *Genes*
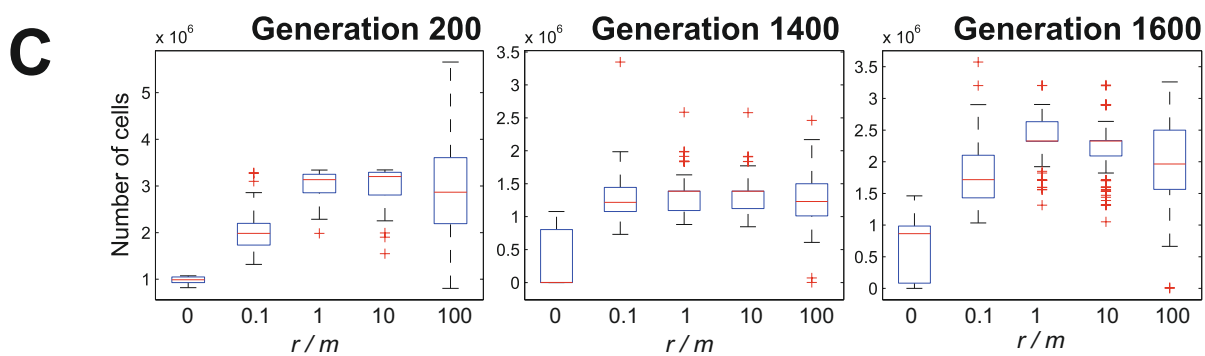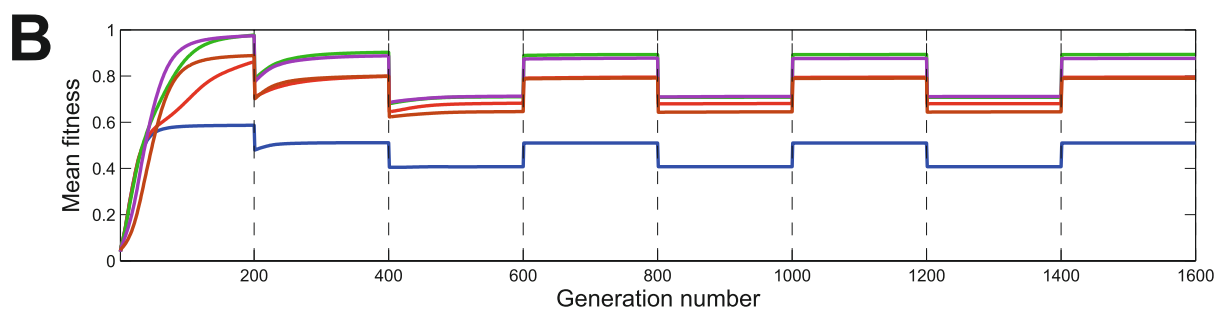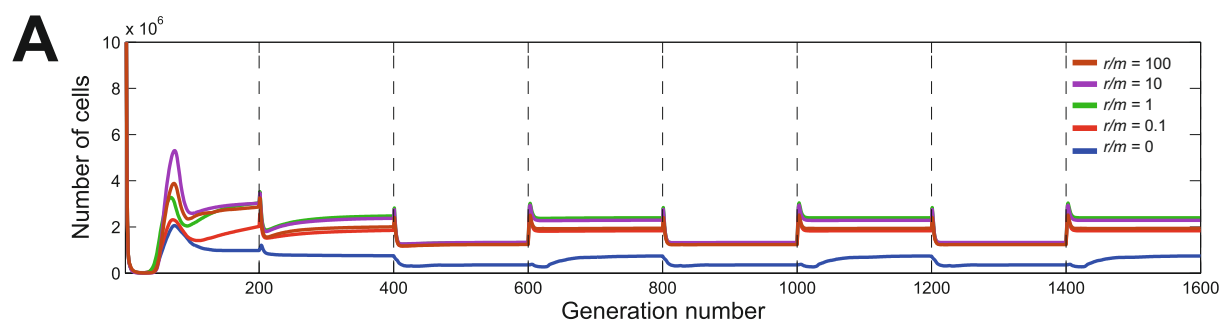686 **3:** 261-277.

688 Sheppard SK, Didelot X, Jolley KA, Darling AE, Pascoe B, Meric G *et al* (2013a).
689 Progressive genome-wide introgression in agricultural Campylobacter coli. *Mol*
690 *Ecol* **22:** 1051-1064.

692 Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ *et al* (2013b).
693 Genome-wide association study identifies vitamin B5 biosynthesis as a host
694 specificity factor in Campylobacter. *Proc Natl Acad Sci U S A* **110:** 11923-11927.
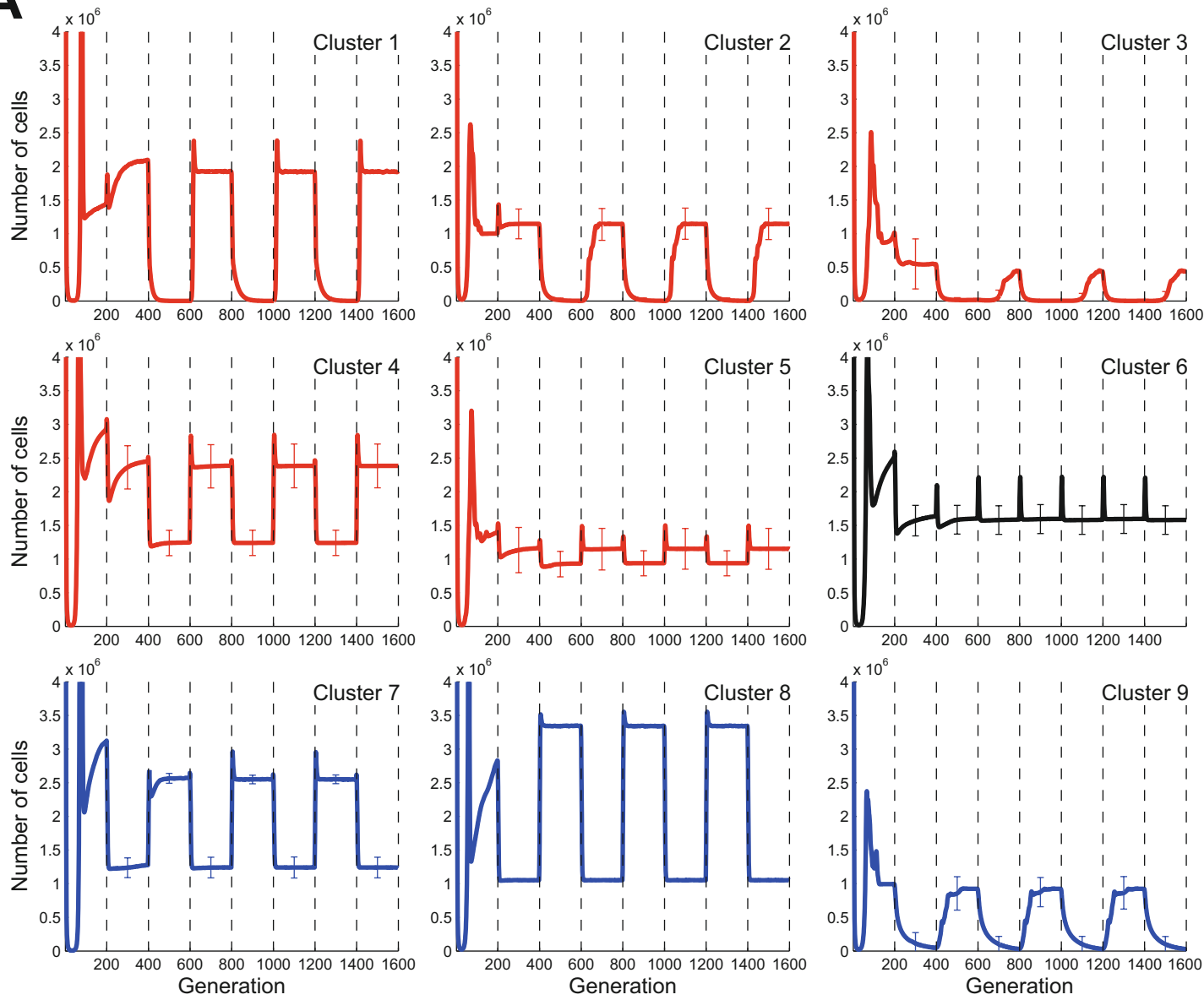
696 Sheppard SK, Cheng L, Meric G, de Haan CP, Llarena AK, Marttinen P *et al*
697 (2014). Cryptic ecology among host generalist Campylobacter jejuni in domestic
698 animals. *Mol Ecol* **23:** 2442-2451.
699

700  Spratt BG, Zhang Q-Y, Jones DM, Hutchison A, Brannigan JA, Dowson CG
701  (1989). Recruitment of a penicillin-binding protien gene from *Neisseria*
702  *flavescens* during the emergence of penicillin resistance in *Neisseria*
703  *meningitidis*. *Proceedings of the National Academy of Sciences USA* **86:** 8988-
704  8992.
705

706  Stireman JO, 3rd (2005). The evolution of generalization? Parasitoid flies and the
707  perils of inferring host range evolution from phylogenies. *J Evol Biol* **18:** 325-336.
708

709  Van Tienderen PH (1991). Evolution of Generalists and Specialist in Spatially
710  Heterogeneous Environments. *Evolution* **45:** 1317-1331.
711

712  Vos M, Didelot X (2009). A comparison of homologous recombination rates in
713  bacteria and archaea. *ISME J* **3:** 199-208.
714

715  Weismann A (1904). *The Evolutionary Theory*, 1 edn. Edward Arnold: London.
716

717  Wiedenbeck J, Cohan FM (2011). Origins of bacterial diversity through horizontal
718  genetic transfer and adaptation to new ecological niches. *FEMS Microbiol Rev*
719  **35:** 957-976.
720

721  Wilson DJ, Gabriel E, Leatherbarrow AJ, Cheesbrough J, Gee S, Bolton E *et al*
722  (2009). Rapid evolution and the importance of recombination to the gastroenteric
723  pathogen Campylobacter jejuni. *Mol Biol Evol* **26:** 385-397.
724

725  Woolhouse ME, Taylor LH, Haydon DT (2001). Population biology of multihost
726  pathogens. *Science* **292:** 1109-1112.
727

728  Zhang H, Gao S, Lercher MJ, Hu S, Chen WH (2012). EvolView, an online tool
729  for visualizing, annotating and managing phylogenetic trees. *Nucleic Acids Res*
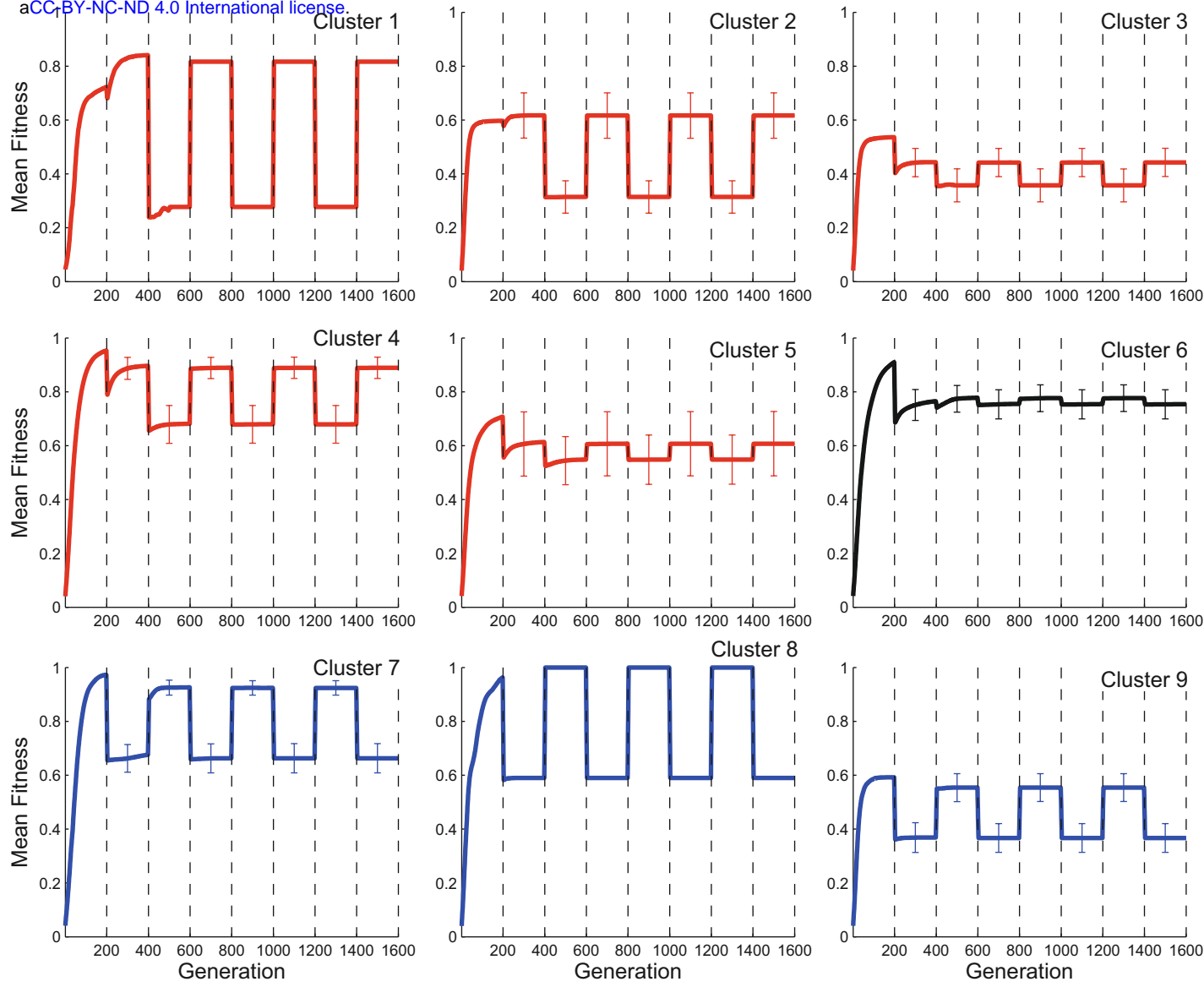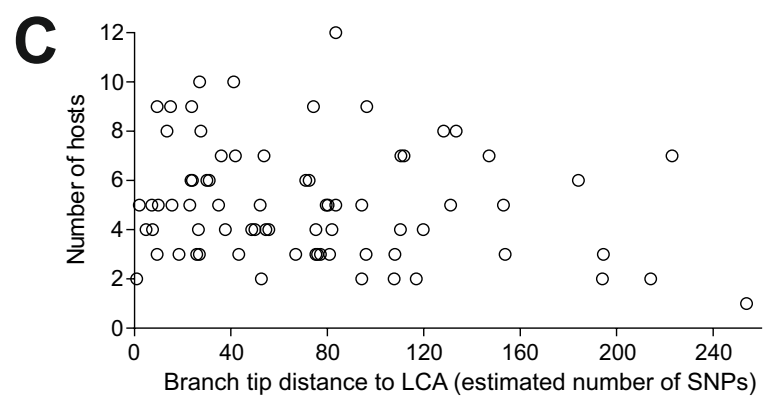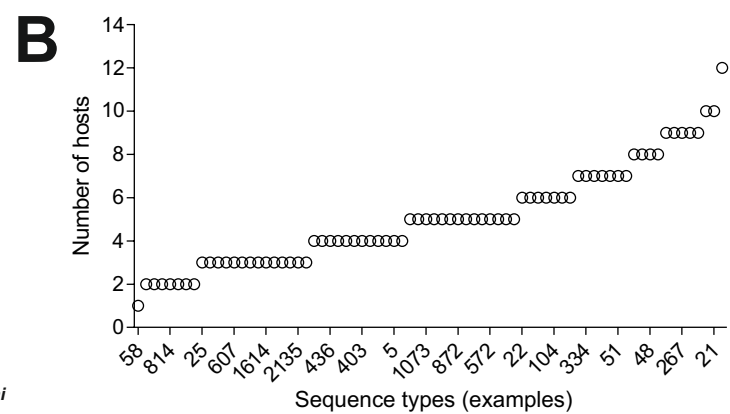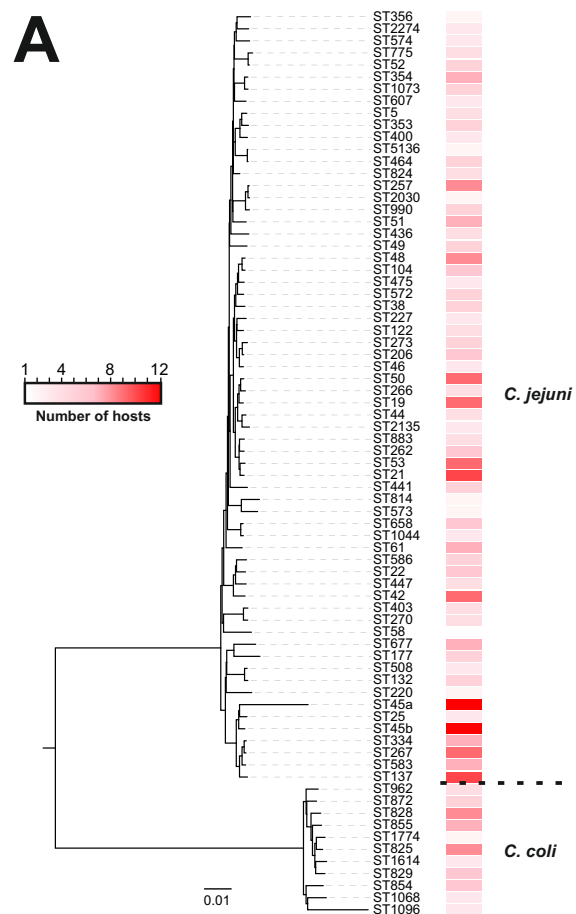730  **40:** W569-572.
731
732

**A**

**Table 1. Proportion of the representative clusters of population profiles at various recombination to mutation ratios**. For each ratio ($r / m$), the percentage of simulations that followed particular representative patterns (clusters 1-9 from Figure 3) from a total of 200 simulations were indicated. Ecological groups were inferred from Figure 3.

| Cluster | Inferred ecological group | Number of simulations (%) | | | | |
|---------|---------------------------|---------|-----------|---------|----------|-----------|
| | | $r/m$=0 | $r/m$=0.1 | $r/m$=1 | $r/m$=10 | $r/m$=100 |
| 1 | Specialist Host 1 | 0 | 0 | 0 | 0 | 0.6 |
| 2 | Specialist Host 1 | 41.1 | 0 | 0 | 0 | 1.2 |
| 3 | Specialist Host 1 | 17.3 | 0 | 0 | 0 | 0.6 |
| 4 | Specialist Host 1 | 0 | 48.2 | 91.8 | 86.4 | 56.7 |
| 5 | Specialist Host 1 | 21.8 | 18.3 | 0 | 1 | 18.7 |
| 6 | Generalist | 0 | 33 | 7.7 | 11.5 | 19.9 |
| 7 | Specialist Host 2 | 0 | 0 | 0.5 | 1 | 0.6 |
| 8 | Specialist Host 2 | 0 | 0.5 | 0 | 0 | 0 |
| 9 | Specialist Host 2 | 19.8 | 0 | 0 | 0 | 1.8 |