Cross-species evolutionary rewiring the enteric bacterium 1 in *Campylobacter* 2

- 3
- Aidan J. Taylor¹§, Guillaume Méric²†§, Koji Yahara³, Ben Pascoe², Leonardos Mageiros^{2,4}, 4
- Evangelos Mourkas², Jessica K Calland², Santeri Puranen⁵, Matthew D. Hitchings⁶, Sion 5
- 6
- Bayliss², Keith A. Jolley⁷, Carolin M. Kobras¹, Nicola J. Williams⁸, Arnoud H. M. van Vliet⁹, Julian Parkhill¹⁰, Martin C. J. Maiden⁷, Jukka Corander^{5,11,12}, Daniel Falush¹³, Xavier 7
- Didelot¹⁴, David J. Kelly^{1*}, Samuel K. Sheppard^{2,7*} 8
- 9 10
- 11 ¹Department of Molecular Biology and Biotechnology, The University of Sheffield, Sheffield S10 12 2TN, United Kingdom;
- ²The Milner Centre for Evolution, Department of Biology and Biochemistry, University of Bath, Bath 13 14 BA2 7AY, UK;
- ³Antimicrobial Resistance Research Centre, National Institute of Infectious Diseases, Tokyo, Japan; 15
- ⁴Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria 16 17 3004. Australia:
- 18 ⁵Department of Mathematics and Statistics, Helsinki Institute for Information Technology, University
- 19 of Helsinki, Helsinki, Finland;
- 20 ⁶Swansea University Medical School, Institute of Life Science, Swansea SA2 8PP, UK;
- 21 ⁷Department of Zoology, University of Oxford, South Parks Road, Oxford OX1 3PS, UK;
- 22 ⁸Department of Epidemiology and Population Health, Institute of Infection & Global Health,
- 23 University of Liverpool, Leahurst Campus, Wirral, UK;
- 24 ⁹School of Veterinary Medicine, University of Surrey, Guildford GU2 7AL, UK;
- 25 ¹⁰Department of Veterinary Medicine, University of Cambridge, Madingley Rd, Cambridge CB3 0ES,
- 26 UK;
- 27 ¹¹Department of Biostatistics, University of Oslo, Oslo, Norway;
- 28 ¹²Parasites and Microbes, Wellcome Sanger Institute, Cambridge, UK;
- ¹³The Centre for Microbes, Development and Health, Institut Pasteur of Shanghai, Shanghai, China; 29
- 30 ¹⁴School of Life Sciences and Department of Statistics, University of Warwick, Coventry CV4 7AL,
- 31

UK.

- 32
- 33
- 34 *Corresponding author(s): David J. Kelly, d.kelly@sheffield.ac.uk; Samuel K. Sheppard, 35 s.k.sheppard@bath.ac.uk;
- 36 †Present address: Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute,
- 37 75 Commercial Rd, Melbourne 3004, Victoria, Australia. 38
- 39 § Contributed equally.
- 40

41 Keywords: comparative genomics, introgression, admixture, epistasis, evolution, Campylobacter, 42 microbiology

- 43
- 44 Running title: Epistasis and recombination in Campylobacter
- 45
- 46

47 Abstract

48

49 The lateral transfer of genes among bacterial strains and species has opposing effects, 50 conferring potentially beneficial adaptations whilst introducing disharmony in coadapted 51 genomes. The prevailing outcome will depend upon the fitness cost of disrupting established 52 epistatic interactions between genes. It is challenging to understand this in nature because it 53 requires population-scale analysis of recombination and genomic coadaptation, and 54 laboratory confirmation of the functional significance of genotype variation. By assigning the 55 ancestry of DNA in the genomes of two species of the world's most common enteric bacterial 56 pathogen, we show that up to 28% of the Campylobacter coli genome has been recently 57 introgressed from Campylobacter jejuni. Then, by quantifying covariation across the genome 58 we show that >83% of putative epistatic links are between introgressed C. jejuni genes in 59 divergent genomic positions (>20kb apart), consistent with independent acquisition. Much of 60 this covariation is between 16 genes, with just 5 genes accounting for 99% of epistatic SNP 61 pairs. Laboratory mutagenesis and complementation cloning assays demonstrated functional 62 links between these genes, specifically related to formate dehydrogenase (FDH) activity. 63 These findings suggest that the genetic confederations that define genomic species may be 64 transient. Even for complex traits such as central metabolism in the bacterial cell, conditions 65 can arise where epistatic genes can be decoupled, transferred and reinstated in a new genetic 66 background. 67

68

69 Introduction

70

71 Complex biological life is made possible by the interaction of genes. In bacteria, a vast array 72 of genes, found together in almost infinite combinations, has allowed colonization of virtually 73 every conceivable habitat on earth. While mutation remains the engine of genetic novelty, for 74 most bacteria adaptation also involves the acquisition of genes from other strains and species 75 through horizontal gene transfer (HGT)¹ potentially conferring new phenotypes to future 76 generations. In some well documented cases, a single nucleotide substitution or acquisition of 77 a small number of genes, can prompt new evolutionary trajectories with striking outcomes 78 such as the emergence of virulent or antibiotic resistant strains². With such dynamic genomic 79 architecture, it may be tempting (and possibly useful³) to consider genes as independent units 80 that 'plug and play' innovation into recipient genomes. This is clearly an oversimplification. 81 In fact, genomes are highly interactive wherein the effect of one gene depends on another 82 (epistasis). Therefore, it is likely that some introduced changes will disrupt gene networks 83 and be costly to the original coadapted genetic background, particularly for complex 84 phenotype innovation involving multiple genes.

85

Understanding how epistasis influences the evolution of phenotype diversity has preoccupied researchers since the origin of population genetics⁴⁻¹⁰, with much emphasis placed upon the 86 87 relative amounts of recombination and epistatic effect sizes 11,12 . In sexual populations, such 88 as outbreeding metazoans, genetic variation is shuffled at each generation so genes can only 89 90 rise in frequency if they have high mean fitness across genetic backgrounds. This means that 91 it is unlikely that multiple distinct epistatic allele combinations will be maintained in the 92 same population and barriers to gene flow, such as geographic isolation, may be required for 93 marked phenotypic diversification⁸. In bacteria, however, rapid clonal reproduction allows 94 multiple independent beneficial allele pairs to rise to high frequency in a single population. 95 For example, in common enteric bacteria such as Escherichia coli, Salmonella enterica and *Campylobacter jejuni* the doubling time in the wild has been estimated at around 24 hours or 96 less^{13,14}. Therefore, though HGT occurs in these organisms¹⁵, even in highly recombinogenic 97 C. $jejuni^{13,16}$, there will likely be many millions of bacterial generations between 98 99 recombination events at a given locus. This allows mutations that are beneficial only in 100 specific genetic backgrounds to establish in a single population and linkage disequilibrium to form between different epistatic pairs¹⁷. 101

102

103 In this coadapted genomic landscape, recombination presents something of a paradox. On the 104 one hand, it promotes adaptation by conferring novel functionality on the recipient genome¹⁸ 105 and reduces competition between clones that carry different beneficial mutations (clonal 106 interference), on the other it potentially introduces disharmonious genes or genecombinations that will be discriminated against by selection¹⁹. Indeed, negative epistatic 107 interactions between genes with different evolutionary histories have been proposed as a barrier to recombination^{4-10,20}, particularly between species. However, interspecies 108 109 recombination is common in bacteria^{18,21}. One of the most conspicuous examples of this can 110 111 be seen in the common animal gut bacterium *Campylobacter*, which is among the most 112 prolific causes of human bacterial gastroenteritis worldwide²². Specifically, introgression 113 between the two most important pathogenic species, C. jejuni and C. coli, has led to the emergence of a globally distributed 'hybrid' C. coli lineage²³ that is responsible for almost all 114 115 livestock and human infections with this species.

116

117 This evolutionary scenario provides an ideal context for addressing the paradox of 118 recombination in coadapted genomes. First, because up to 23% of the core genome of

119 common *C. coli* has been recently introgressed from *C. jejuni*²⁴, potentially disrupting 120 epistatic interactions. Second, because *C. jejuni* and *C. coli* have undergone an extended 121 period of independent evolution (85% average nucleotide identity), therefore recombined 122 sequence is conspicuous in the genome. Third, because the outcompetition of ancestral strains 123 by introgressed *C. coli* is inconsistent with traditional assumptions of selection against hybrid 124 lineages.

125

126 Here, we analyse the genomes of isolates representing known diversity within C. jejuni and 127 C. coli and perform chromosome painting analysis to assign ancestry across the genomes of 128 introgressed C. coli. We then perform a systematic screen for long-range covariation in the 129 core genome to identify SNPs that are in strong linkage disequilibrium in independent genetic 130 backgrounds and, therefore, are not explained by the clonal frame of the population. Finally, 131 we investigate the function of covarying genes that account for most of the interactions and 132 confirm epistasis among coadapted genes. By combining these analyses, we demonstrate how 133 multiple interacting genes from one species can transfer to another, modifying the co-adapted 134 genomic landscape of the recipient and promote a natural admixture-mediated genetic 135 revolution^{19,25}.

- 136
- 137

138 **Results**

139

140 Campylobacter populations are highly structured with intermediate sequence clusters

141 Pan-genome analysis identified the presence and variation of every automatically annotated 142 gene from every genome. This revealed a core genome of 631 gene orthologues in all 143 Campylobacter isolates in this study. There were 1287 genes common to all C. jejuni isolates 144 and 895, 1021 and 1272 common to C. coli clades 1, 2 and 3, respectively. Consistent with 145 previous studies²⁴, neighbour-joining and ClonalFrameML trees based on genes within 146 concatenated core genome alignments revealed population structure in which C. *jejuni* and C. 147 *coli* clade 2 and 3 isolates each formed discrete clusters (Figure 1A, Supplementary Figure 1). 148 Isolates designated as C. coli clade 1 were found in three clusters on the phylogeny: 149 unintrogressed ancestral strains, and the ST-828 and ST-1150 clonal complexes which 150 account for the great majority of strains found in agriculture and human disease²⁴. Pan-151 genome analysis quantified the increase in unique gene discovery as the number of sampled 152 genomes increased (Figure 1B). For all sequence clusters there was evidence of an open pan-153 genome with a trend towards continued rapid gene discovery within sequence clusters with 154 fewer isolates. There was considerable accessory genome variation between species and 155 clades, potentially associated with important adaptive traits (Supplementary Figure 2) and 156 there was evidence that the average number of genes per genome was greater in C. coli ST-157 828 complex isolates than in *C. jejuni* (Figure 1C).

158

159 There is substantial introgression within the *C. coli* core genome

160 The large genetic distances among C. coli clade 1 isolates (Figure 1A), have been shown to 161 be a consequence of the import of DNA from C. jejuni rather than accumulation of mutations 162 during a prolonged period of separate evolution²³. Using chromosome painting to infer the 163 coancestry of core-genome haplotype data from CC-828 and CC-1150 isolates gave a 164 detailed representation of the recombination-derived chunks from each C. jejuni donor group 165 to each recipient individual (Supplementary Figure 3). The majority of introgressed SNPs 166 were rare, occurring in fewer than 50 recipient genomes (Figure 2A). However, a large 167 proportion of the introgressed C. coli clade 1 genomes contained DNA of C. jejuni ancestry in >98% of recipient isolates. Consistent with previous estimates²⁴, these regions where 168

introgressed DNA was largely fixed within the *C. coli* population, occurred across the
genome and comprised up to 15% and 28% of the CC-828 and CC-1150 isolate genomes
respectively (Figure 2A). When considering donor groups, the majority of introgressed DNA
in *C. coli* involves genes that are present in multiple *C. jejuni* lineages (core genes) (Figure
2B).

174

175 Having identified C. jejuni ancestry within C. coli genomes, we investigated the sequence of 176 events responsible for introgression. Most introgressed SNPs are found at low frequency in 177 both clonal complexes. However, there was evidence of SNPs that are introgressed in both 178 complexes as well as high frequency lineage specific introgression in both CC-828 and CC-179 1150 (Figure 2C). Specifically, 25% of the C. jejuni DNA found in >98% of CC-828 isolates 180 was also found in CC-1150 (Figure 2D), implying that this genetic material was imported by 181 the common ancestor(s) of both complexes. Subsequent to the divergence of these two 182 complexes, introgression continued with nearly 75% of C. jejuni DNA present in one 183 recipient clonal complex and not the other. This is consistent with an evolutionary history in 184 which there was a period of progressive species and clade divergence reaching approximately 185 12% at the nucleotide level between C. jejuni and C. coli and around 4% between the 186 three C. coli lineages. More recently, changes to the patterns of gene flow led one 187 C. coli clade 1 lineage to import substantial quantities of C. jejuni DNA, and further lineage-188 specific introgression gave rise to two clonal complexes (CC-828 and CC-1150) that 189 continued to accumulate C. jejuni DNA, independently creating the population structure 190 observed today (Figure 2E).

191

192 The high magnitude introgression into C. coli clade 1 isolates has introduced thousands of 193 nucleotide changes to the core genome. However, divergence in bacteria may be uneven 194 across the genome. First, because recombination is more likely to occur in regions where donor and recipient genomes have high nucleotide similarity^{21,26,27}. Second, because of 195 'fragmented speciation'²⁸, in which gene flow varies in different parts of the genome, such as 196 197 regions responsible for adaptive divergence, leading to phylogenetic incongruence among genes. Consistent with previous estimates²⁴, we found that the three *C. coli* clades had similar 198 199 high divergence with C. jejuni across the genome, ranging from 68% to 98% nucleotide 200 identity for individual genes (Figure 2F), implying a period of divergence with low levels of 201 gene flow. We found no evidence that high genetic differentiation between the species 202 prevented recombination. While there was some evidence that more recombination occurred 203 in regions of low nucleotide divergence (between unintrogressed C. coli clade 1 and C. 204 *jejuni*), introgression occurred across the genome at sites with varying levels of nucleotide 205 identity (Figure 2F). This level of recombination has greatly increased overall genetic 206 diversity across the genome in C. coli clade 1 and introduced changes that have potential 207 functional significance.

208

209 Much of the putative epistasis occurs between SNPs in introgressed genes

210 ClonalFrameML analysis revealed the importance of homologous recombination in 211 generating sequence variation within the introgressed C. coli. Estimates of the relative 212 frequency of recombination versus mutation ($R/\theta=0.43$), mean recombination event length 213 $(\delta = 152 \text{ bp})$ and average amount of polymorphism per site in recombined fragments (v=0.07), 214 imply that recombination has had an effect (r/m) 4.57 times higher than *de novo* mutation during the diversification of CC-828. This is consistent with previous analysis and confirmed 215 recombination as the major driver of molecular evolution in C. coli^{13,24,29}. The continuous 216 217 time Markov chain model for the joint evolution of pairs of biallelic sites on a phylogenetic 218 tree (Supplementary Figure 4) was applied to investigate patterns of covariance for all pairs

of sites >20kb apart (Figure 3A). For most biallelic sites there were few branches on the tree where substitutions occurred, so that their evolution is compatible with separate evolution on the same clonal frame. However, 2874 covarying pairs evolved more frequently together than would be expected (*p*-value 10^{-8}) if they had evolved independently based on the tree, and hence indicated patterns of putative epistasis.

224

225 Among them, the location of 2618 putative epistatic pairs of sites was compared to the 226 inferred ancestry (unintrogressed C. coli or C. jejuni) of sequence across the genome of CC-227 828 and CC-1150 C. coli strains (Figure 3B, Supplementary Data 1). For each epistatic pair, 228 the major and minor haplotype were defined if there was haplotype polymorphism between 229 C. jejuni and CC-828 and CC-1150 C. coli. This allowed quantification of the number of 230 covarying sites that occurred between an ancestral C. coli (unintrogressed) and an 231 introgressed C. jejuni allele, two introgressed alleles, and sites that do not segregate by 232 species. Strikingly, the breakdown of the major and minor haplotype combinations among the 233 2618 epistatic pairs (Figure 3C, Source Data) shows the major haplotype for 83.5% of 234 putative epistatic SNP pairs was C. jejuni indicating that both co-varying sites had C. jejuni 235 ancestry, consistent with epistasis between introgressed ancestral C. jejuni sequence at 236 divergent genomic positions. Investigation of the genes containing co-varying sites revealed 237 that 2187 SNP pairs were in 16 genes with just five genes accounting for 99.1% of them 238 (Figure 3D, Supplementary Data 2, Supplementary Figure 5).

239

240 Genomic context and physiological role of epistatically linked genes

241 The five genes accounting for the majority of epistatic interactions (cj1167, cj1168c, 242 cj1171c/ppi, cj1507c/modE and cj1508c/fdhD) were investigated for their physiological role 243 in C. jejuni. FdhD and ModE are proteins involved in the biogenesis of formate 244 dehydrogenase (FDH). The FDH complex (FdhABC) oxidises formate to bicarbonate to 245 generate electrons that fuel cellular respiration. Formate is an abundant electron donor produced by host microbiota and an important energy source for Campylobacter in vivo^{30,31}. 246 247 The remaining three genes, *cj1167* (annotated incorrectly as *ldh*, lactate dehydrogenase), 248 cj1168c and cj1171c (ppi) are also grouped together on the genome, where cj1167 and 249 cj1168c are adjacent but with the open reading frames (ORFs) convergent and overlapping, 250 while *ppi* is upstream, separated by two non-epistatically linked genes (c_{j1169c} and c_{j1170c} , 251 Figure 4A). Considering the genomic arrangement, it is therefore clear that the putative 252 epistatic links uncovered in this study essentially occur across two loci in the genome 253 (*fdhD/modE* and *cj1167/cj1168c/ppi*), with each of the latter three genes linked with both 254 fdhD and modE (Figure 4A). Given the known function of fdhD and modE in biogenesis of 255 the FDH complex, we hypothesised that cj1167/cj1168c/ppi might also have some role in 256 FDH biogenesis or activity in order to form a functional epistatic connection. We therefore 257 constructed deletion mutants to investigate the possible role of these genes in FDH activity.

258

259 Initially, each of the mutants and their parental wildtype (C. jejuni NCTC11168) were grown 260 in rich media (Muller-Hinton broth) and their formate dependent oxygen consumption rates 261 determined (Figure 4B). cj1167, cj1168c and ppi mutants demonstrated wildtype levels of 262 FDH activity, while activity in both *fdhD* and *modE* mutants was abolished. In order to 263 confirm that the phenotype of the *fdhD* and *modE* mutants was not due to a polar effect on the 264 surrounding *fdh* locus, these mutants were genetically complemented by reintroduction of a 265 second copy of the wildtype gene into the rRNA locus, which restored near wildtype levels of 266 FDH activity in both cases (Figure 4B).

267

268 As neither *cj1167*, *cj1168c* or *ppi* mutants showed altered FDH activity in cells grown in rich 269 media, we considered that their function may be related to an FDH-specific nutrient 270 requirement as would likely be found in vivo. Since the formate oxidising subunit of FDH, 271 FdhA, specifically requires a molybdo- or tungsto-pterin (Mo/W) cofactor and a selenocysteine (SeC) residue for catalysis³², Mo, W or Se supply presented possible targets. 272 273 cj1168c encodes a DedA family integral membrane protein of unknown function. DedA 274 proteins are solute transporters widespread in bacteria but are mostly uncharacterised³³. 275 However, a homologue of c_{j1168c} in the heavy metal specialist beta-proteobacterium 276 *Cupriavidus metallidurans* has been shown to be involved in selenite (SeO₃²⁻) uptake³⁴. We 277 therefore speculated that Cj1168 could be a selenium oxyanion transporter that supplies Se 278 for SeC biosynthesis. To test this, FDH activities were measured in *cill68c* mutant and 279 parental wildtype strains grown in minimal media with limiting concentrations of selenite or 280 selenate (SeO₄²⁻). The data in Figure 4C shows that the cj1168c mutant displayed 281 significantly reduced FDH activity after growth with selenite in the low nM range, and this 282 phenotype was partially restored by genetic complementation. We therefore designated 283 cj1168c as selF (selenium transporter for formate dehydrogenase). However, although this 284 phenotype does suggest that SelF is a selenium importer, another unrelated selenium 285 transporter, FdhT (Cj1500), has previously been documented in C. jejuni³⁵, which is not 286 epistatic with *fdhD* or *modE*. In contrast to this previous report we found considerable 287 residual FDH activity still remained in an fdhT deletion mutant, which was fully restored to 288 wildtype levels by complementation (Figure 4D).

289

290 Finally, we tested whether the residual FDH activity in our *fdhT* mutant was due to selenium 291 uptake by SelF. An *fdhT selF* double mutant was generated and assayed for FDH activity 292 after growth in minimal media containing limiting concentrations of selenite or selenate 293 (Figure 4D). The *fdhT selF* double mutant demonstrated a significant additional reduction in 294 FDH activity over the fdhT single mutant, a phenotype that was partially restored by 295 complementing the double mutant with *selF*. Complementation of the double mutant with 296 *fdhT* returned FDH activity to near wildtype levels (Figure 4D). Taken together, our data 297 suggests that both FdhT and SelF facilitate selenium acquisition in C. jejuni, possibly 298 representing low and high affinity transporters, respectively (Figure 4E). 299

300 Discussion

301

302 Hybridization between distantly related organisms can bring together new gene combinations and traits that potentially allow adaptation in a single evolutionary leap³⁶. However, 303 304 introducing disharmony among genes that have coevolved in epistasis can lead to reduced hybrid fitness, limiting the chance of this type of evolution by saltation^{8,10}. Evidence for these 305 306 contrasting paradigms comes largely from eukaryotes, but among prokaryotes, widespread 307 HGT may seem to contradict assumptions about selection against hybrid lineages. For 308 example, the transfer of mobile resistance genes between bacterial strains and species can 309 confer a clear adaptive benefit, in a single step, allowing the recipient lineages to proliferate 310 in the presence of antibiotics. Clearly, therefore, epistasis is not an absolute barrier to long-311 range recombination in bacteria.

312

There is ongoing debate about the extent to which the galaxy of accessory genes that are variously present or absent in many bacterial genomes constitute a cache of mobile genes from which innovation can be drawn and transferred between strains or species³⁷⁻⁴⁰. However, considering HGT of accessory genes does not fully address the extent to which recombination is constrained by epistatic fitness interactions between genes. First, because

genes or gene clusters introduced as autonomous elements, encoding specific novel traits, may cause minimal disruption to other essential cellular functions. Second, even when introduced genes result in a concomitant change elsewhere in the genome⁴¹, the fitness cost may be outweighed by the benefit in a given niche. Third, and most importantly, because most recombination in bacteria is between homologous sequences²¹, where a given gene is replaced by another version of the same gene.

324

325 *Campylobacter* is an ideal model organism for considering barriers to gene flow as 326 interspecies recombination is well documented within the core genome. Consistent with previous studies 23,24,42 , we found that a large proportion (15-28%) of the core genome of the 327 328 two common C. coli lineages found in livestock and human disease (CC-828 and CC-1150) 329 originated in C. jejuni. It is challenging to explain this level of genome-wide gene flow 330 between species after a prolonged period of diversification (~15% nucleotide divergence) as 331 it would likely disrupt coadapted gene networks. The evident success of hybrid C. coli in the 332 agricultural niche suggests that the accumulation of C. jejuni DNA was not detrimental, but 333 to specifically address if genome coadaptation was a barrier to recombination requires 334 quantification of covariation among alleles.

335

336 Most bacteria, including *Campylobacter*, have highly genetically structured populations reflecting both neutral and adaptive evolutionary processes¹⁸. Therefore, a phylogenetically 337 338 naive statistical model of coadaptation would afford equal significance to linkage 339 disequilibrium resulting from selection and common ancestry. Even in recombinogenic 340 bacteria such as C. coli, HGT is not sufficiently common to scramble the genome and abolish 341 non-random SNP associations resulting from clonal population structure. Consistent with 342 other models⁴³, the statistical test developed here accounts for the amount of covariation that 343 would be expected based upon the clonal frame and identifies the same combinations of 344 alleles in independent genetic backgrounds, thus providing evidence for coadaptation.

345

346 By combining the results from the chromosome painting and covariation models we were 347 able to quantify the frequency of C. coli and C. jejuni SNPs in covarying allele pairs across 348 the genome of introgressed C. coli. Recent interspecies admixture results in allele pairs, that 349 correspond to C. jejuni – C. coli (and vice versa) and C. jejuni – C. jejuni. Comparing 350 haplotype frequency provides the opportunity to contrast the disruptive and beneficial effects 351 of homologous recombination in the bacterial genome. Under a neutral model, where 352 recombination does not disrupt beneficial epistatic interactions, there would be more C. jejuni 353 -C. coli than C. jejuni -C. jejuni covarying allele pairs in the recipient genome. Clearly, 354 relatively low levels of covariation between ancestral C. coli and introgressed C. jejuni SNPs 355 is expected, consistent with selection against disharmonious gene combinations in the 356 coadapted recipient genome. However, the finding that C. jejuni - C. jejuni allele pairs 357 constituted >83% of covarying introgressed haplotypes was striking. It is possible that in 358 some cases both sites were introgressed in a single recombination event as bacteria can import very large pieces of DNA (>100kb)⁴⁴⁻⁴⁶ but in *Campylobacter* LD for pairs of sites 359 decreases with distance to approximately 20kbp and then remains at the same level for very 360 361 distant sites⁴⁷. Therefore, the divergent genome position of allele pairs (>20kb) implies that 362 they were acquired independently. It follows, therefore, that the acquisition of the first 363 introgression event was not fatal to the recipient genome, and was either mildly detrimental, 364 neutral or beneficial. Acquisition of the second member of the pair then potentially enhanced 365 the fitness restoring the integrated C. *jejuni* – C. *jejuni* coevolutionary unit. 366

367 Investigating covariation provides clues about genome plasticity and the potential benefit of 368 introgression for niche adaptation. However, to directly address if genes are in epistasis it is 369 necessary to test their function and confirm coadaptation. We demonstrated a functional link 370 amongst the top scoring co-varying gene pairs with FDH, a key enzyme allowing the utilisation of formate as an electron donor in vivo^{30,31}. In this study, FdhD and ModE were 371 372 shown to be essential for FDH activity. While FdhD is a sulfur-transferase known to be required for the insertion of the pterin cofactor into FdhA⁴⁸ (Figure 4E), ModE is a 373 374 transcriptional repressor that has been shown previously only to regulate the Mo/W uptake genes mod and $tup^{49,50}$. However, the unexpected abolished FDH activity in a modE mutant 375 376 indicates further functions for ModE in FDH biogenesis which warrant future investigation. 377 Searching for functional links between *fdhD/modE* and *cj1167/selF/ppi*, revealed that a *selF* 378 mutant strain had significantly reduced FDH activity under conditions of selenite limitation, a 379 phenotype consistent with SelF being a Se oxyanion transporter and that functionally links 380 SelF with FdhD/ModE. We suggest that selF rather than fdhT is epistatic because SelF 381 confers an additional benefit for SeC biosynthesis (essential for FDH activity) under 382 conditions of selenium limitation, for example as may be found in the host (Figure 4E). 383 cj1167 encodes a cytoplasmic NADPH dependent 2-oxoacid dehydrogenase but the current 384 genome database annotation as lactate dehydrogenase (Ldh) is incorrect and its function is 385 unknown⁵¹. There is no precedent for the involvement of such an enzyme in bacterial FDH or 386 SeC biogenesis and we obtained no evidence for a functional connection between Cj1167 and 387 FDH activity. However, the overlapping convergent gene arrangement of cj1167 and cj1168c 388 (selF) suggests a transcriptional architecture that might dictate these genes both form similar 389 epistatic dependencies even if Cj1167 is not required for FDH activity. Finally, cj1171c (ppi) 390 encodes a cytoplasmic peptidyl-prolyl *cis-trans* isomerase of the cyclophilin family. PPIases 391 are general protein folding catalysts that often have pleiotropic and redundant functions⁵² and 392 we note that our C. jejuni ppi deletion mutant showed no growth defect as well as no 393 reduction in FDH activity. It is possible that if Cj1171 does help promote the folding of e.g. 394 FdhD or ModE, analysis of a simple deletion mutant may not reveal this if another PPIase 395 can substitute in that genetic background.

396

397 Explaining the aberrant genome architecture among introgressed C. coli is challenging. 398 Understanding the selective value of genes that promote proliferation depends on the overall genetic environment¹⁹. Our findings are consistent with an evolutionary scenario where an 399 400 ancestral C. coli lineage underwent a niche transition and the surviving lineages (CC-828 and 401 CC-1150) gained access to C. jejuni DNA (Figure 3E&F). As the adaptive landscape of the 402 genome changed, potentially decoupling epistatic interactions that were previously selected, 403 new gene combinations could be introduced by recombination and tested in the C. coli 404 genetic background. Recent studies have shown that this type of genetic rewiring may be more common than previously thought⁵³⁻⁵⁵ and when it leads to the proliferation of a hybrid 405 organism it can be associated with the colonization of a new niche. Intensive livestock 406 407 systems represent a possible explanation for the genetic revolution in C. coli. Host ecology 408 can dramatically affect the evolution of gut-dwelling organisms, including Campylobacter⁵⁶ 409 and colonization of this niche could promote the conditions necessary to promote 410 introgressed C. coli. Whether host ecology is a factor or not, it is clear that conditions can 411 arise where coadapted genes in a highly interactive bacterial genome can be transferred 412 between species and reinstated as a single evolutionary unit in a new genome. This suggests 413 that epistasis is not an absolute barrier to genome-wide recombination in structured bacterial 414 populations.

- 415
- 416 Methods

417

418 Isolates, genome sequencing and assembly

419 A total of 973 isolates were used in this study, 827 from C. coli and a selection of 146 from a 420 diversity of C. jejuni clonal complexes (Supplementary Data 3). Isolates were sampled 421 mostly in the United Kingdom to maximise environmental and riparian reservoirs and thus 422 the representation of genetic diversity in C. coli. Isolates were stored in a 20% (v/v) glycerol 423 medium mix at -80°C and subcultured onto *Campylobacter* selective blood-free agar 424 (mCCDA, CM0739, Oxoid). Plates were incubated at 42°C for 48 h under microaerobic 425 conditions (5% CO2, 5% O2) generated using a CampyGen (CN0025, Oxoid) sachet in a 426 sealed container. Subsequent phenotype assays were performed on Brucella agar (CM0271, 427 Oxoid). Colonies were picked onto fresh plates and genomic DNA extraction was carried out 428 using the QIAamp® DNA Mini Kit (QIAGEN; cat. number: 51306) according to the 429 manufacturer's instructions. DNA was eluted in 100–200 µl of the supplied buffer and stored 430 at -20° C. DNA was quantified using a Nanodrop spectrophotometer and high-throughput 431 genome sequencing was performed on a MiSeq (Illumina, San Diego, CA, USA), using the 432 Nextera XT Library Preparation Kit with standard protocols involving fragmentation of 2 μ g 433 genomic DNA by acoustic shearing to enrich for 600 bp fragments, A-tailing, adapter ligation and an overlap extension PCR using the Illumina 3 primer set to introduce specific tag 434 435 sequences between the sequencing and flow cell binding sites of the Illumina adapter. DNA 436 cleanup was carried out after each step to remove $DNA \square < \square 150$ bp using a 1:1 ratio of 437 AMPure[®] paramagnetic beads (Beckman Coulter, Inc., USA). Short read paired-end data was assembled using the *de novo* assembly algorithm, SPAdes (version 3.10.0)⁵⁷. All novel 438 439 genome sequences (n=475) generated for use in this study are available on NCBI BioProjects 440 PRJNA689604 and PRJEB11972. These were augmented with 498 previously published genomes and accession numbers for all genomes can be found in Supplementary Data 3^{16,24,29,47,56,58-61} 441 442

443

444

445 Genome archiving, pan-genome content analyses and phylogenetic reconstruction

446 Contiguous genome sequence assemblies were individually archived on the web-based database platform BIGSDB⁶² and sequence type (ST) and clonal complex (CC) designation 447 448 were assigned based upon the C. *jejuni* and C. *coli* multi-locus sequence typing scheme⁶³. To 449 examine the full pan-genome content of the dataset, a reference pan-genome list was assembled as previously described⁶⁴. Briefly, genome assemblies from all 973 genomes in 450 this study were automatically annotated using the RAST/SEED platform⁶⁵, the BLAST 451 452 algorithm was used to determine whether coding sequences from this list were allelic variants 453 of one another or 'unique' genes, with two alleles of the same gene being defined as sharing 454 >70% sequence identity on >10% of the sequence length. The prevalence of each gene in the 455 collection of 973 genomes was determined using BLAST with a positive hit in a genome 456 being defined as a local alignment of the reference sequence with the genomic sequence of 457 >70% identity on >50% of the length, as previously described⁶⁶. The resulting matrix was 458 analysed for differentiating core and accessory genome variation. Genes present in all 459 genomes were concatenated to produce a core-genome alignment, used for subsequent phylogenetic reconstructions. Phylogenetic trees were reconstructed using an approximation 460 of maximum-likelihood phylogenetics in FastTree2⁶⁷. This tree was used as an input for 461 ClonalFrameML⁶⁸ to produce core genome phylogenies with branch lengths corrected for 462 463 recombination.

464

465 Inference of introgression

466 All 973 genomes were aligned to a full reference sequence of C. coli strain CVM29710. We conducted imputation for polymorphic sites with missing frequency $\leq 10\%$ using BEAGLE⁶⁹ 467 as previously reported⁷⁰. A total of 286,393 gapless SNPs (~17% of the average C. coli 468 469 genome size) were used for recombination analyses. The coancestry of genome-wide haplotype data was inferred based on alignments using chromosome painting 470 and FineStructure⁷¹ as previously described⁷². Briefly, ChromoPainter was used to infer 471 472 chunks of DNA donated from a list of 33 donor groups normalised for sample size to each of 473 677 ST-CC-828 and 12 CC-1150 recipient haplotypes. Results were summarised into a 474 coancestry matrix containing the number of recombination-derived chunks from each donor 475 to each recipient individual. FineStructure was then used for 100,000 iterations of both the 476 burn-in and Markov chain Monte Carlo chain to cluster individuals based on the co-ancestry 477 matrix. The results are visualized as a heat map with each cell indicating the proportion of 478 DNA "chunks" (a series of SNPs with the same expected donor) a recipient receives from 479 each donor.

480

481 Analysis of covariation in bacterial genomes

482 Non-random allele associations can result from selection and clonal population structure. To 483 control for the latter, our approach identified SNP combinations in independent genetic 484 backgrounds by accounting for the sequence variation associated with the inferred phylogeny. Based on the alignment of 677 genomes of C. coli CC-828, a first phylogenetic tree was 485 created using PhyML⁷³. ClonalFrameML⁶⁸ was then applied to correct the tree by accounting 486 487 for the effect of recombination, and also to infer the ancestral sequence of each node. 488 Covariance was assessed for pairs of biallelic sites across the genome using a Continuous 489 Time Markov Chain (CTMC) model as follows. Briefly, let A and a denote the two alleles of 490 the first site and B and b denote the two alleles of the second site, so that there are four states 491 in total for the pair of sites (ab, Ab, aB and AB). The four substitution rates from A to a, from 492 a to A, from B to b and from b to B are not assumed to be identical, to allow for differences in 493 substitution rates in different parts of the genome and also to allow for non-equal rates of 494 forward and backward substitution (for example as a result of recombination opportunities). 495 Assuming no epistatic effect between the two sites (ε =1), the model M₀ has four free 496 parameters (α_1 , α_2 , β_1 and β_2) representing independent substitutions at the two sites. We 497 expand model M_0 with an additional fifth parameter $\varepsilon > 1$ into model M_1 which is such that the 498 state AB where the first site is allele A and the second site is allele B is favored relative to the 499 other three sites ab, aB and Ab. Specifically, the state AB has a probability increased by a 500 factor ε^2 in the stationary distribution of the CTMC of model M₁ compared to model M₀.

501

502 Both models M_0 (with 4 parameters) and M_1 (with 5 parameters) are fitted to the data using 503 maximum likelihood techniques, where the likelihood is equal to the product for every branch 504 of the tree of the state at the bottom of the branch given the state at the top. The two fitted 505 models M_0 and M_1 are then compared using a likelihood-ratio test (LRT) as follows: since M_0 506 is nested with M_1 , two times their difference in log-likelihood is expected to be distributed 507 according to a chi-square distribution with number of degrees of freedom equal to the 508 difference in their dimensionality, which is one. This LRT returns a p-value for the 509 significance of a covariation effect, and a Bonferroni correction is applied to determine a 510 conservative cutoff of significance that accounts for multiple testing. Furthermore, the test is 511 applied only to pairs of sites separated by >20kb to reduce the chance that they were the result of a single recombination event, consistent with estimates of the length of recombined 512 DNA sequence in quantitative bacterial transformation experiments^{44,74} and evidence from 513 514 Campylobacter genome analyses that show that LD for pairs of sites decreases with distance to approximately 20kbp and then remains at the same level for very distant sites⁴⁷. It is still 515

516 possible of course that rare recombination events would stretch $20kbp^{44-46}$, but for this to 517 have an effect on the analysis of epistasis it would have to have happened several times for 518 the same pairs of sites against different genomic background which becomes quite unlikely 519 just by chance. This phylogenetically aware approach to testing for covariance presents the 520 advantage to naturally account for both population structure and the effect of 521 recombination⁷⁵. The script implementing this coevolution test is available in R at: 522 https://github.com/xavierdidelot/campy.

523

524 Quantifying covariation between recombined and unrecombined genomic regions

525 The results of the introgression and covariation analyses were combined so that for each pair 526 of significantly covarying SNPs (*p*-value $<10^{-8}$), haplotype frequency was calculated among 527 the 689 recipient introgressed C. coli clade-1 strains as well as among the donor C. coli 528 (ancestral) and C. jejuni strains, respectively. If the most frequent haplotype of the pair is the 529 same between the donor C. coli (ancestral) and C. jejuni, it was classified as 'no 530 polymorphism'. Otherwise, if the most frequent haplotype accounted for >90% among the 531 recipients, it was classified as either 'C. jejuni (>90%)' or 'C. coli (>90%)' if it was the same 532 as that of donor C. jejuni or C. coli (ancestral) (inset in Figure 3C). If the most frequent 533 haplotype accounted for $\leq 90\%$ among the recipients, the top two most frequent haplotypes 534 (written as major and minor haplotype in this manuscript) were indicated as either "C. jejuni / 535 C. coli", "C. jejuni / other", "C. coli / C. jejuni", "C. coli / other", "other / C. coli", "other / C. 536 *jejuni*", and "other / other", and the frequency of the major and minor haplotypes were 537 calculated. For example, where the haplotype frequencies were as follows, AA=285, 538 TA=192, TG=181, AG=27, A=2, --=1, -A=1, AA is the major haplotype, frequency of 539 which is 41.3%

540

541 Mutagenesis and complementation cloning

542 Genes cj1167, cj1168c (here designated selF for selenium transport for formate 543 dehydrogenase), cj1171c (ppi), cj1507c (modE), cj1508c (fdhD) and cj1500 (fdhT) were 544 deleted by allelic exchange mutagenesis, with the majority of the open reading frame 545 replaced by an antibiotic resistance cassette. Mutagenesis plasmids were generated by the 546 isothermal assembly method using the HiFi system (NEB, UK). In brief, flanking regions of 547 target genes were PCR amplified from genomic DNA using primers with adaptors 548 homologous to either the backbone vector pGEM3ZF or the antibiotic resistance cassette 549 (Supplementary Data 4). pGEM3ZF was linearised by digestion with HincII. The kanamycin 550 and chloramphenicol resistance cassettes were PCR amplified from pJMK30 and pAV35, 551 respectively⁷⁶. Four fragments consisting of linearised pGEM3ZF, antibiotic resistance 552 cassette and 2 flanking regions were combined in equimolar amounts and mixed with 2 x 553 HiFi reagent (NEB, UK) and incubated at 50°C for 1 hour. The fragments combine such that 554 the gene fragments flank the antibiotic resistance cassette, in the same transcriptional 555 orientation, within the vector. Mutagenesis plasmids were transformed into C. jejuni NCTC 556 11168 by electroporation. Spontaneous double-crossover recombinants were selected for 557 using the appropriate antibiotic and correct insertion into the target gene confirmed by PCR 558 screening. For genetic complementation of mutants, genes cj1168c (selF), cj1507c (modE), 559 cj1508c (fdhD) and cj1500 (fdhT) were PCR amplified from genomic DNA, restriction digested with MfeI and XbaI, then ligated into similarly digested pRRA⁷⁷ (Supplementary 560 561 Data 4). The orientation of insertion allowed the target gene to be expressed constitutively by 562 a chloramphenicol resistance gene-derived promoter within the vector. Complementation 563 plasmids were transformed into C. jejuni by electroporation. Spontaneous double-crossover 564 recombinants were selected for using apramycin and correct insertion into the ribosomal 565 locus confirmed by PCR screening.

566

567 Growth of C. jejuni

Microaerobic growth cabinets (Don Whitley, UK) were maintained at 42°C with an 568 569 atmosphere of 10% O₂, 5% CO₂ and 85% N₂ (v/v). C. jejuni was grown on Columbia-base agar containing 5% v/v defibrinated horse blood. Selective antibiotics were added to plates as 570 appropriate at the following concentrations: 50 μ g ml⁻¹ kanamycin, 20 μ g ml⁻¹ 571 chloramphenicol, 60 µg ml⁻¹ apramycin. Muller-Hinton (MH) broth supplemented with 20 572 mM L-serine was used as a rich medium. Minimal medium was prepared from a supplied 573 574 MEM base (51200-38, Thermo Scientific, UK) with the following additions: 20 mM L-575 serine, 0.5 mM sodium pyruvate, 50 µM sodium metabisulfite, 4 mM L-cysteine. HCl, 2 mM 576 L-methionine, 5 mM L-glutamine, 50 µM ferrous sulfate, 100 µM ascorbic acid, 1 µM 577 vitamin B12, 5 µM sodium molybdate, 1 µM sodium tungstate. Selenium was then added as 578 appropriate from stocks of sodium selenate or sodium selenite prepared in dH₂O. For assays, 579 cells were washed and suspended in sterile phosphate-buffered saline (PBS, Sigma-Aldrich).

580

581 **Respiration rates with formate**

582 Cells were first grown in MH broth for 12 hours, then washed thoroughly in PBS before 583 inoculating minimal media without an added selenium source. The appropriate concentrations 584 were determined by serial dilution trials and it was subsequently found that C. jejuni has a 585 strong preference for selenite over selenate, as equivalent FDH activity requires some 1000-586 fold greater concentration of selenate than selenite (Supplementary Figure 6). These cultures 587 were grown for 8 hours before the cells were thoroughly washed again, then used to inoculate 588 further minimal media, with a selenium source added as appropriate, and grown for 10 hours. 589 This passaging was necessary to remove all traces of selenium from the inoculum, such that 590 control cultures without selenium added had negligible (FDH) activity. Assay cultures were 591 again thoroughly washed before the equivalent of 20 ml at an optical density of 0.8 at 600 nm 592 was finally suspended in 1 ml of PBS. Formate-dependent oxygen consumption by whole 593 cells was measured in a Clark-type electrode using 20 mM sodium formate as electron 594 donor. The electrode was calibrated with air-saturated PBS assuming 220 nmol dissolved O_2 595 ml⁻¹ at 37 °C. In the electrode, 200 μ l of the dense cell suspension was added to 800 μ l air-596 saturated PBS for a final volume of 1 ml. The chamber was sealed and the suspension 597 allowed to equilibrate for 2 minutes. The assay was initiated by the addition of 20 μ l of 1 M 598 sodium formate (prepared in PBS) and the rate of oxygen consumption recorded for 90 s. The 599 total protein concentration of the cell suspensions was determined by Lowry assay and the 600 specific rate of formate-dependent oxygen consumption expressed as nmol oxygen consumed 601 \min^{-1} mg⁻¹ total protein.

602

603 Data availability

604 Short-read sequence data for all isolates sequenced in this study are deposited in the sequence 605 read archive (SRA) and can be found associated with NCBI BioProjects PRJNA689604 606 (https://www.ncbi.nlm.nih.gov/bioproject/PRJNA689604) and **PRJEB11972** 607 (https://www.ncbi.nlm.nih.gov/bioproject/PRJEB11972). These were augmented with 498 608 previously published genomes and assembled genomes are available on Figshare 609 (doi.org/10.6084/m9.figshare.13521683). Accession numbers for all genomes are included in 610 Supplementary Data 3. Source data are provided for this paper.

611

612 Acknowledgements

613 This work was supported by Wellcome Trust grants 088786/C/09/Z and Medical Research 614 Council (MRC) grants MR/M501608/1 and MR/L015080/1 awarded to SKS. Bacterial 615 sampling was funded by Food Standards Agency project FS101087. AJT was supported by a 616 Biotechnology and Biological Sciences Research Council grant (BB/S014497/1) awarded to 617 DJK. GM was supported by a NISCHR Health Research Fellowship (HF-14-13). JC was 618 supported by the ERC grant no. 742158. Computational calculations were performed using 619 computer servers from the MRC CLIMB project (UK), at National Institute of Genetics 620 (Japan) and the Human Genome Center of the Institute of Medical Science (University of 621 Tokyo, Japan) and at HPC Wales (UK).

622

623 Author Contributions

S.S., X.D. D.F. and K.Y. conceived and designed the study. S.S., N.W., A.V. and M.M.
collected samples. G.M., A.T., L.M., M.H. and B.P. carried out Laboratory work. B.P., M.H.,
K.J., M.M., J.P. and S.S. supported data archiving. G.M., A.T., X.D., K.Y., L.M., S.P., S.S.
and J.C. analysed the data. M.M., J.P., J.C., D.K. and D.F. contributed to data interpretation.
S.S., B.P., G.M., C. K., A.T. and D.K. wrote the paper.

629

630 Competing Interests

- 631 The authors declare no competing interests.
- 632
- 633
- 634

635 Figure legends

636

637 Figure 1. Population genomics of C. jejuni and C. coli. (A) Phylogenetic tree reconstructed 638 using neighbour-joining on a whole-genome alignment of 973 C. jejuni and C. coli isolates. 639 Introgressed C. coli clades are represented with red (CC-828, n=677) and purple (CC-1150, 640 n=12) circles, unintrogressed clade 1 (n=35) is shown in pink, clade 2 (n=45) in yellow and 641 clade 3 (n=58) in green. A set of 146 C. jejuni genomes belonging to 30 clonal complexes (4 642 to 5 isolates per ST) are shown in blue. Recipient and donor populations, used to infer 643 introgression in chromosome painting analysis, are indicated. The scale bar represents the 644 number of substitutions per site. (B) Accumulation curves of the clade-specific pan-genome 645 content of Campylobacter lineages, using the same colour code as panel A. Randomized 646 genome sampling was used to obtain the average number of genes for each comparison (plain 647 lines) and standard deviations (dotted lines). (C) The average number of genes/genome, 648 identified by BLAST, is shown for C. jejuni and the different C. coli clades. Asterisks denote 649 significant difference between distributions as inferred from a Dunn's multiple comparisons 650 test after a Kruskal-Wallis test, as follows: ***: p<0.001; ****: p<0.0001.

651

652 Figure 2. Genome-wide introgression from C. jejuni to C. coli. (A) Summary of 653 introgressed C. jejuni SNPs in C. coli CC-828 (n=677, red) and CC-1150 (n=12, purple) 654 genomes using ChromoPainter; the number of introgressed core SNPs (coloured histograms; 655 left y-axis) and core genes (white histograms; right y-axis) for a range of recipient strains 656 proportions (at least 1, more than 50% and more than 98%) is shown. (B) The number of 657 genes with different frequencies of maximum SNP introgression/gene in C. coli as a function 658 of gene frequency in C. jejuni. Highly introgressed genes in CC-828 and CC-1150 tend to be 659 core in C. jejuni. (C) Density plot (n=1000 bins) of specific and shared introgression events 660 in CC-828 (x-axis) and CC-1150 (y-axis). The frequency of SNP introgression/gene is shown 661 for both lineages. Close blue lines denote a high density of points. (D) Shared introgression 662 between C. coli CC-828 and CC-1150. The number of SNPs being shared between the two 663 lineages at various frequencies is shown in y-axis. (E) Diagram of *Campylobacter* species 664 and clade (C1*,C2, C3) divergence with arrows indicating introgression from C. jejuni into 665 C. coli (i) clade 1, (ii, iii) CC-828 and CC-1150, (iv, v) subsequent clonal expansion and 666 ongoing introgression. (F) Pairwise nucleotide identity between C. jejuni and ancestral 667 (unintrogressed) clade 1 C. coli core genes (black circles). Genes found to be introgressed in 668 clade 1 CC-828 are highlighted in blue.

669

670 Figure 3. Covariation in introgressed C. coli genomes. (A) CC-828 and CC-1150 C. coli 671 genomes were analysed using a continuous time Markov chain (CTMC) model and 672 covariation was assessed for pairs of biallelic sites separated by at least 20kb along the 673 genome while accounting for the effect of population structure and recombination. There are 674 many biallelic sites that do not change often on the tree and few that do. Putative epistatic 675 sites change more frequently than average with biallelic pairs found together on multiple 676 branches. (B) Miami plot of each polymorphic site showing the maximum p-value for 677 covarying biallelic pairs (>20kb apart) and the frequency of introgression in CC-828 and CC-678 1150. (C) The frequency of major and minor haplotype combinations (inset) among the 2578 679 pairs of covarying SNPs in the 689 C. coli clade-1 recipient genomes, revealing that the 680 majority of long range covariation was between introgressed C. jejuni sites. (D) The position 681 of putative epistatic sites mapped on the C. coli CVM29710 reference for covarying C. 682 *jejuni-C. jejuni* SNPs (red) in 16 gene pairs (a to l), and other haplotype combinations (grey). 683 (E and F) An evolutionary scenario for the observed patterns of covariation and introgression 684 in natural C. coli populations: (i) C. jejuni (blue) and unintrogressed C. coli (red) co-exist

with genomes (internal circles) harbouring haplotype pairs (x-x) that segregate by species; (ii) Horizontal gene transfer, HGT, occurs (R1) disrupting covarying genetic elements and reducing the relative fitness of introgressed *C. coli* to varying degrees (grey arrow), few strains retain mixed *C. coli* - *C. jejuni* haplotypes; (iii) HGT continues (R2) and, where recombined mixed haplotypes survived, ancestral *C. jejuni* haplotype pairs are reinstated in introgressed *C. coli*; (iv,v) introgressed *C. coli* outcompete unintrogressed strains.

691

692 Figure 4. Genomic context and physiological roles of introgressed epistatically linked 693 genes. (A) Genome organisation and percentage of co-varying SNP pairs (internal legend). 694 (B-D) FDH activity of whole cells determined by oxygen consumption rates in a Clark-type 695 electrode (nmol oxygen consumed per minute per mg of total protein) for (B) cells grown in 696 rich media (excess selenium), (C) cells grown in minimal media with 0.5 nM sodium selenite, 697 and (D) cells grown in minimal media with either 5 nM sodium selenite (left, open bars) or 5 698 μ M sodium selenate (right, hashed bars). All data are means of at least 4 independent 699 determinations and error bars are SD. *** denotes p value of <0.001 by students t-test. (E) 700 Model for epistatically linked genes involving FDH biogenesis and activity. Host derived 701 formate is converted to bicarbonate in the periplasm by the FDH complex to release electrons 702 which are transferred from the iron sulfur (FeS) cluster of FdhA to the *b*-type hemes of FdhC 703 and into the menaquinone (MK) pool where they can ultimately be used to reduce molecular 704 oxygen via terminal oxidases. FdhA contains a selenocysteine residue and Mo/W-pterin 705 cofactor (W/MoCo) at its active site, both of which are essential for catalysis. ModE is a 706 DNA binding regulator which represses the Mo and W transporters Mod and Tup to regulate 707 the cellular pool of Mo/W. W/MoCo is generated by the Moa pathway and inserted into apo-708 FdhA by the sulphur-transferase FdhD. Environmental selenite (the most abundant oxyanion) 709 or selenate diffuses into the periplasm where it must be actively imported to the cytoplasm by 710 FdhT and SelF. Putatively, FdhT is low affinity and functions efficiently when ample 711 selenium is available. When selenium is limited, SelF can import sufficient selenium to 712 maintain FDH production and activity. Cytoplasmic selenium is converted to 713 selenophosphate by SelD, which is used by SelA to generate tRNA-SeC from tRNA-Ser. 714 During translation of FdhA, tRNA-SeC is incorporated by the specific elongation factor SelB. 715 Apo-FdhA with W/MoCo and SeC inserted is then transported to the periplasm by the TAT translocation system and incorporated into the FDH complex (See ³¹ for a review of FDH in 716 717 C. jejuni).

- 718
- 719

720

721

722	Refere	ences
723	1	Ochman, H., Lawrence, J. G. & Groisman, E. A. Lateral gene transfer and the nature
724		of bacterial innovation. <i>Nature</i> 405 , 299-304 (2000).
725	2	Koonin, E. V., Makarova, K. S. & Aravind, L. Horizontal gene transfer in
726		prokaryotes: quantification and classification. Annu. Rev. Microbiol. 55, 709-742
727		(2001).
728	3	Haldane, J. B. A defense of beanbag genetics. <i>Perspect Biol Med</i> 7, 343-359 (1964).
729	4	Bateson, W. in Darwin And Modern Science 85–101 (Cambridge University Press,
730		1909).
731	5	Fisher, R. A. The correlation between relatives on the supposition of Mendelian
732		inheritance. Trans. R. Soc. Edinb 52, 399–433 (1918).
733	6	Wright, S. Evolution in Mendelian populations. <i>Genetics</i> 16, 97-159 (1931).
734	7	Dobzhansky, T. Studies on hybrid sterility. II. Localization of sterility factors in
735		Drosophila Pseudoobscura hybrids. Genetics 21, 113-135 (1936).
736	8	Dobzhansky, T. Genetics And The Origin Of Species. (Columbia University Press,
737		1937).
738	9	Muller, H. J. Bearing of the Drosophila work on systematics. The new systematics,
739		185-268 (1940).
740	10	Muller, H. Isolating mechanisms, evolution, and temperature. Biol. Symp. 6, 71-125
741		(1942).
742	11	Kimura, M. Attainment of quasi linkage equilibrium when gene frequencies are
743		changing by natural selection. Genetics 52, 875-890 (1965).
744	12	Neher, R. A. & Shraiman, B. I. Competition between recombination and epistasis can
745		cause a transition from allele to genotype selection. Proc. Natl. Acad. Sci. 106, 6866-
746		6871 (2009).
747	13	Wilson, D. J. et al. Rapid evolution and the importance of recombination to the
748		gastroenteric pathogen Campylobacter jejuni. Mol. Biol. Evol. 26, 385-397 (2009).
749	14	Gibson, B., Wilson, D. J., Feil, E. & Eyre-Walker, A. The distribution of bacterial
750		doubling times in the wild. Proc. R. Soc. B. 285 (2018).
751	15	Vos, M. & Didelot, X. A comparison of homologous recombination rates in bacteria
752		and archaea. <i>Isme J</i> 3 , 199-208 (2009).
753	16	Calland, J. K. et al. Quantifying bacterial evolution in the wild: a birthday problem for
754		Campylobacter lineages. Preprint at
755		https://www.biorxiv.org/content/10.1101/2020.12.02.407999v1 (2020).
756	17	Arnold, B. J. et al. Weak epistasis may drive adaptation in recombining bacteria.
757		<i>Genetics</i> 208 , 1247-1260 (2018).
758	18	Sheppard, S. K., Guttman, D. S. & Fitzgerald, J. R. Population genomics of bacterial
759		host adaptation. Nat. Rev. Genet. 19, 549-565 (2018).
760	19	Mayr, E. in Evolution as a Process 157–180 (Allen & Unwin, 1954).
761	20	Orr, H. A. Dobzhansky, Bateson, and the genetics of speciation. Genetics 144, 1331-
762		1335 (1996).
763	21	Fraser, C., Hanage, W. P. & Spratt, B. G. Recombination and the nature of bacterial
764		speciation. Science 315, 476-480 (2007).
765	22	Sheppard, S. K. et al. Campylobacter genotyping to determine the source of human
766		infection. Clin. Infect. Dis. 48, 1072-1078 (2009).
767	23	Sheppard, S. K., McCarthy, N. D., Falush, D. & Maiden, M. C. Convergence of
768		Campylobacter species: implications for bacterial evolution. Science 320, 237-239
769		(2008).
770	24	Sheppard, S. K. et al. Progressive genome-wide introgression in agricultural
771		Campylobacter coli. Mol. Ecol. 22, 1051-1064 (2013).

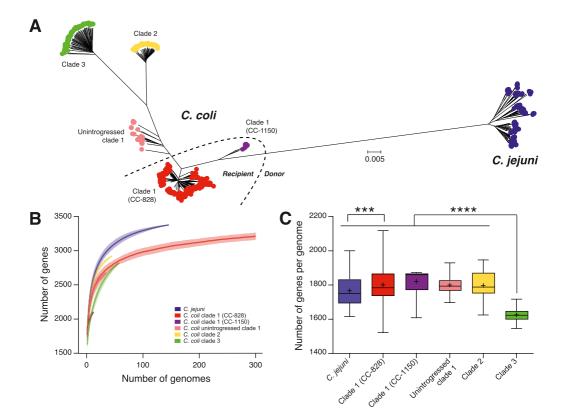
- Provine, W. B. in *Genetics, Speciation And The Founder Principle* 43-76 (Oxford University Press, 1989).
- Cohan, F. M. in *Genetics of Mate Choice: From Sexual Selection to Sexual Isolation*359-370 (Springer, 2002).
- Ansari, M. A. & Didelot, X. Inference of the properties of the recombination process
 from whole bacterial genomes. *Genetics* 196, 253-265 (2014).
- Retchless, A. C. & Lawrence, J. G. Phylogenetic incongruence arising from
 fragmented speciation in enteric bacteria. *Proc. Natl. Acad. Sci.* 107, 11453-11458
 (2010).
- Dearlove, B. L. *et al.* Rapid host switching in generalist *Campylobacter* strains erodes
 the signal for tracing human infections. *Isme J* 10, 721-729 (2016).
- Weerakoon, D. R., Borden, N. J., Goodson, C. M., Grimes, J. & Olson, J. W. The role
 of respiratory donor enzymes in *Campylobacter jejuni* host colonization and
 physiology. *Microb. Pathog.* 47, 8-15 (2009).
- Taylor, A. J. & Kelly, D. J. The function, biogenesis and regulation of the electron
 transport chains in *Campylobacter jejuni*: New insights into the bioenergetics of a
 major food-borne pathogen. *Adv. Microb. Physiol.* **74**, 239-329 (2019).
- Smart, J. P., Cliff, M. J. & Kelly, D. J. A role for tungsten in the biology of *Campylobacter jejuni*: tungstate stimulates formate dehydrogenase activity and is transported via an ultra-high affinity ABC system distinct from the molybdate transporter. *Mol. Microbiol.* **74**, 742-757 (2009).
- Doerrler, W. T., Sikdar, R., Kumar, S. & Boughner, L. A. New functions for the ancient DedA membrane protein family. *J. Bacteriol.* 195, 3-11 (2013).
- Ledgham, F., Quest, B., Vallaeys, T., Mergeay, M. & Covès, J. A probable link
 between the DedA protein and resistance to selenite. *Res. Microbiol.* 156, 367-374
 (2005).
- Shaw, F. L. *et al.* Selenium-dependent biogenesis of formate dehydrogenase in *Campylobacter jejuni* is controlled by the *fdhTU* accessory genes. *J. Bacteriol.* 194, 3814-3823 (2012).
- 801
 36
 Rieseberg, L. H. Hybrid origins of plant species. Annu. Rev. Ecol. Syst. 28, 359-389

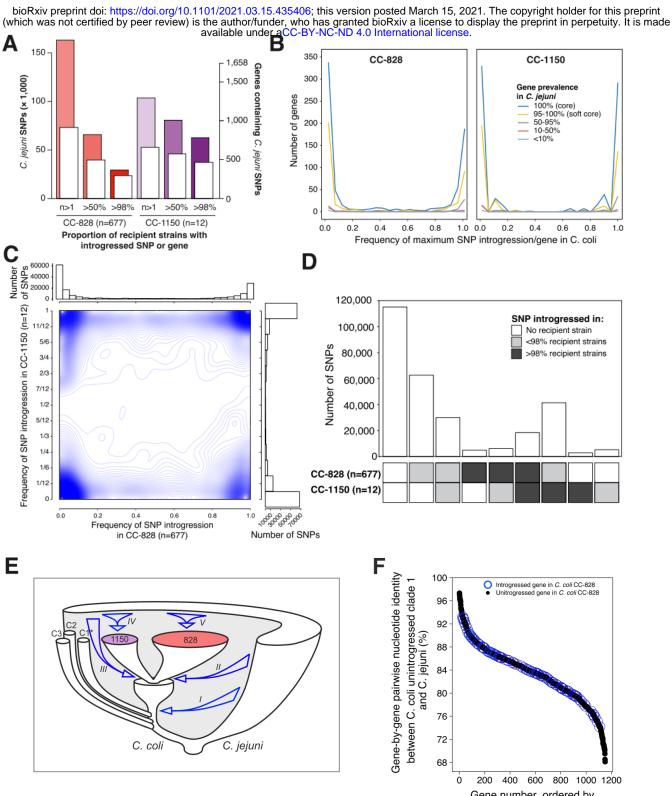
 802
 (1997).
- 80337Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent on804effective population size. *Isme J* 11, 1719-1721 (2017).
- 805 38 Vos, M. & Eyre-Walker, A. Are pangenomes adaptive or not? *Nature Microbiol.* 2, 1576-1576 (2017).
- Shapiro, B. J. The population genetics of pangenomes. *Nature Microbiol.* 2, 15741574 (2017).
- McInerney, J. O., McNally, A. & O'Connell, M. J. Why prokaryotes have
 pangenomes. *Nature Microbiol.* 2, 17040 (2017).
- 41 Levin, B. R., Perrot, V. & Walker, N. Compensatory mutations, antibiotic resistance
 and the population genetics of adaptive evolution in bacteria. *Genetics* 154, 985-997
 (2000).
- Sheppard, S. K., McCarthy, N. D., Jolley, K. A. & Maiden, M. C. J. Introgression in the genus *Campylobacter*: generation and spread of mosaic alleles. *Microbiology* 157, 1066-1074 (2011).
- 817 43 Pensar, J. *et al.* Genome-wide epistasis and co-selection study using mutual
 818 information. *Nucleic Acids Res.* 47, e112-e112 (2019).
- K. K., Barquist, L., Parkhill, J. & Bentley, S. D. A high-resolution view of genome-wide pneumococcal transformation. *PLoS Path.* 8, e1002745 (2012).

- 45 He, M. *et al.* Evolutionary dynamics of *Clostridium difficile* over short and long time
 823 scales. *Proc. Natl. Acad. Sci.* 107, 7527-7532 (2010).
- Power, P. M., Bentley, S. D., Parkhill, J., Moxon, E. R. & Hood, D. W. Investigations
 into genome diversity of *Haemophilus influenzae* using whole genome sequencing of
 clinical isolates and laboratory transformants. *BMC Microbiol.* 12, 273 (2012).
- Yahara, K. *et al.* Genome-wide association of functional traits linked with *Campylobacter jejuni* survival from farm to fork. *Environ. Microbiol.* 19, 361-380
 (2017).
- 48 Arnoux, P. *et al.* Sulphur shuttling across a chaperone during molybdenum cofactor
 maturation. *Nat. Commun.* 6, 6148 (2015).
- Taveirne, M. E., Sikes, M. L. & Olson, J. W. Molybdenum and tungsten in *Campylobacter jejuni*: their physiological role and identification of separate transporters regulated by a single ModE-like protein. *Mol. Microbiol.* **74**, 758-771 (2009).
- Aguilar-Barajas, E., Díaz-Pérez, C., Ramírez-Díaz, M. I., Riveros-Rosas, H. &
 Cervantes, C. Bacterial transport of sulfate, molybdate, and related oxyanions. *BioMetals* 24, 687-707 (2011).
- Thomas, M. T. *et al.* Two respiratory enzyme systems in *Campylobacter jejuni* NCTC
 11168 contribute to growth on l-lactate. *Environ. Microbiol.* 13, 48-61 (2011).
- 52 Ünal, C. M. & Steinert, M. Microbial peptidyl-prolyl cis/trans isomerases (PPIases):
 virulence factors and potential alternative drug targets. *Microbiol. Mol. Biol. Rev.* 78,
 544-571 (2014).
- Solution Structure
 Solution Structure</l
- 84654Arnold, B. *et al.* Fine-scale haplotype structure reveals strong signatures of positive847selection in a recombining bacterial pathogen. *Mol. Biol. Evol.* **37**, 417-428 (2019).
- Wadsworth, C. B., Arnold, B. J., Sater, M. R. A. & Grad, Y. H. Azithromycin
 resistance through interspecific acquisition of an epistasis-dependent efflux pump
 component and transcriptional regulator in *Neisseria gonorrhoeae. mBio* 9, e0141901418 (2018).
- 85256Mourkas, E. et al. Agricultural intensification and the evolution of host specialism in
the enteric pathogen Campylobacter jejuni. Proc. Natl. Acad. Sci. 117, 11018-11028
(2020).854(2020).
- 855 57 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications
 856 to single-cell sequencing. *Journal of computational biology : a journal of*857 *computational molecular cell biology* 19, 455-477 (2012).
- Sheppard, S. K. *et al.* Genome-wide association study identifies vitamin B5
 biosynthesis as a host specificity factor in *Campylobacter. Proc. Natl. Acad. Sci.* 110, 11923-11927 (2013).
- Sheppard, S. K. *et al.* Cryptic ecology among host generalist *Campylobacter jejuni* in domestic animals. *Mol. Ecol.* 23, 2442-2451 (2014).
- 863 60 Pascoe, B. *et al.* Local genes for local bacteria: evidence of allopatry in the genomes
 864 of transatlantic *Campylobacter* populations. *Mol. Ecol.* 26, 4497-4508 (2017).
- 865 61 Pascoe, B. *et al.* Genomic epidemiology of *Campylobacter jejuni* associated with
 866 asymptomatic pediatric infection in the Peruvian Amazon. *PLoS Negl. Trop. Dis.* 14,
 867 e0008533 (2020).
- Bingle, K. E. *et al.* Multilocus sequence typing system for *Campylobacter jejuni*. J. *Clin. Microbiol.* 39, 14-23 (2001).

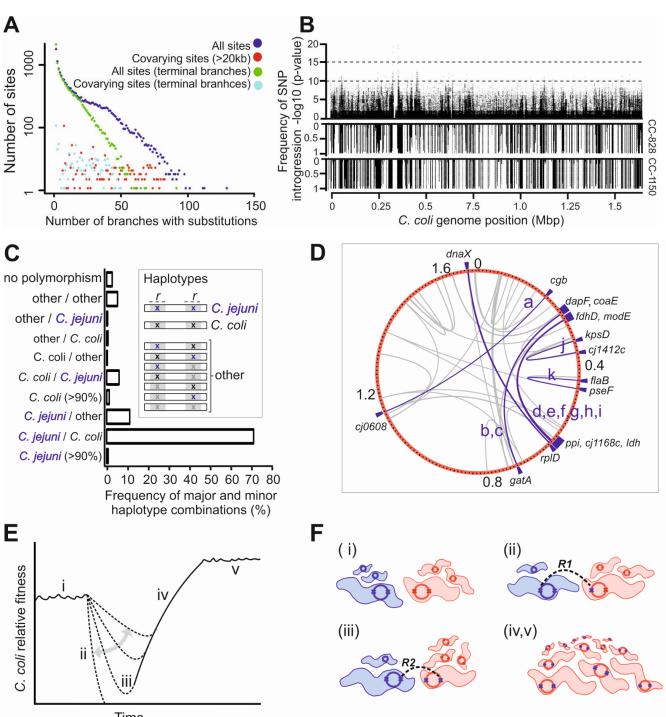
- Méric, G. *et al.* A reference pan-genome approach to comparative bacterial genomics:
 identification of novel epidemiological markers in pathogenic *Campylobacter*. *PLOS One* 9, e92798 (2014).
- Aziz, R. K. *et al.* The RAST server: rapid annotations using subsystems technology. *BMC Genomics* 9, 75 (2008).
- Sheppard, S. K., Jolley, K. A. & Maiden, M. C. J. A gene-by-gene approach to
 bacterial population genomics: whole genome MLST of *Campylobacter. Genes* 3,
 261-277 (2012).
- Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 approximately maximumlikelihood trees for large alignments. *PLOS One* 5, e9490 (2010).
- Bidelot, X. & Wilson, D. J. ClonalFrameML: efficient inference of recombination in
 whole bacterial genomes. *PLoS Comp. Biol.* 11, e1004041 (2015).
- Browning, B. L. & Browning, S. R. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* 84, 210-223 (2009).
- Yahara, K. *et al.* Chromosome painting *in silico* in a bacterial species reveals fine
 population structure. *Mol. Biol. Evol.* 30, 1454-1464 (2013).
- Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population
 structure using dense haplotype data. *PLoS Genet.* 8, e1002453 (2012).
- Yahara, K., Didelot, X., Ansari, M. A., Sheppard, S. K. & Falush, D. Efficient inference of recombination hot regions in bacterial genomes. *Mol. Biol. Evol.* 31, 1593-1605 (2014).
- Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood
 phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* 59, 307-321 (2010).
- Mell, J. C., Shumilina, S., Hall, I. M. & Redfield, R. J. Transformation of natural genetic variation into *Haemophilus influenzae* genomes. *PLoS Path.* 7, e1002151-e1002151 (2011).
- Collins, C. & Didelot, X. A phylogenetic method to perform genome-wide association studies in microbes that accounts for population structure and recombination. *PLoS Comp. Biol.* 14, e1005958 (2018).
- van Vliet, A. H., Wooldridge, K. G. & Ketley, J. M. Iron-responsive gene regulation
 in a *Campylobacter jejuni fur* mutant. *J. Bacteriol.* 180, 5291-5298 (1998).
- 90477Cameron, A. & Gaynor, E. C. Hygromycin B and apramycin antibiotic resistance905cassettes for use in *Campylobacter jejuni*. *PLOS One* **9**, e95084 (2014).

906





Gene number, ordered by pairwise nucleotide identity



Time

