

1 **Plasmids do not consistently stabilize cooperation across**  
2 **bacteria, but may promote pathogen host-range**

3 Anna E. Dewar<sup>1,a,\*</sup>, Joshua L. Thomas<sup>1,a</sup>, Thomas W. Scott<sup>1</sup>, Geoff Wild<sup>2</sup>,  
4 Ashleigh S. Griffin<sup>1</sup>, Stuart A. West<sup>1,b</sup>, Melanie Ghoul<sup>1,b</sup>

5 <sup>1</sup>Department of Zoology, University of Oxford, Oxford, OX1 3SZ, United Kingdom

6 <sup>2</sup>Department of Applied Mathematics, University of Western Ontario, London, Ontario N6A  
7 3K7, Canada

8 a Joint first author

9 b Joint last author

10 \*Corresponding author

11

12 Key words/phrases: extracellular proteins, genetic architecture, horizontal gene transfer,  
13 inclusive fitness, kin selection, secretome.

14

15 **Abstract**

16 Horizontal gene transfer via plasmids could favour cooperation in bacteria, because transfer of  
17 a cooperative gene turns non-cooperative cheats into cooperators. This hypothesis has received  
18 support from theoretical, genomic and experimental analyses. In contrast, with a comparative  
19 analysis across 51 diverse species, we found that genes for extracellular proteins, which are  
20 likely to act as cooperative ‘public goods’, were not more likely to be carried on either: (i)  
21 plasmids compared to chromosomes; or (ii) plasmids that transfer at higher rates. Our results  
22 were supported by theoretical modelling which showed that while horizontal gene transfer can  
23 help cooperative genes initially invade a population, it has less influence on the longer-term  
24 maintenance of cooperation. Instead, we found that genes for extracellular proteins were more  
25 likely to be on plasmids when they coded for pathogenic virulence traits, in pathogenic bacteria  
26 with a broad host-range.

27

28 **Introduction**

29 The growth and success of many bacterial populations depends upon the production of  
30 cooperative ‘public goods’<sup>1-4</sup>. Public goods are molecules whose secretion provides a benefit  
31 to the local group of cells. Examples include iron-scavenging siderophores<sup>5</sup>, exotoxins that

32 disintegrate host cell membranes<sup>6,7</sup>, and elastases that break down connective tissues<sup>8-10</sup>. A  
33 problem is that cooperation can be exploited by ‘cheats’: cells which avoid the cost of  
34 producing public goods but can still use and benefit from those produced by cooperative  
35 cells<sup>3,11,12</sup>. What prevents cheats from outcompeting cooperators, and ultimately destabilising  
36 cooperation?

37

38 In bacteria, some genetic elements are able to move between cells<sup>13</sup>. This horizontal gene  
39 transfer has been suggested as a mechanism to help stabilize the production of cooperative  
40 public goods<sup>14-18</sup> (Figure 1a). If a gene coding for the production of a public good can be  
41 transferred horizontally, it would allow cheats to be ‘infected’ with the cooperative gene and  
42 turned into cooperators. Theoretical models have shown that this can facilitate the invasion of  
43 cooperative genes, in conditions where they would not be favoured on chromosomes<sup>14-18</sup>.  
44 Experiments on a synthetic *Escherichia coli* system have shown that location on a plasmid  
45 helped the gene for a cooperative public good to invade, particularly in structured  
46 populations<sup>18</sup>. In addition, bioinformatic analyses across a range of species found that genes  
47 that code for extracellular proteins, many of which act as public goods, are more likely to be  
48 found on plasmids than the chromosome<sup>15,19,20</sup>.

49

50 There are, however, three potential problems for the hypothesis that horizontal gene transfer  
51 favours cooperation. First, previous bioinformatic analyses made important first steps, but are  
52 not conclusive. One study examined only a single species, which may not be representative of  
53 all bacteria<sup>15</sup>. Two additional studies examined multiple species, but assumed that genes and  
54 genomes from the same and different species can be treated as independent data points, in a  
55 way that could have led to spurious results<sup>19,20</sup>. Statistical tests typically assume that data points  
56 are independent, and even slight non-independence can lead to heavily biased results (type I  
57 errors)<sup>21,22</sup>. There is an extensive literature in the field of evolutionary biology showing that  
58 species share characteristics inherited through common descent, rather than through  
59 independent evolution, and so cannot be considered independent data points<sup>23-25</sup>. Genomes are  
60 nested within species, and genes are nested within genomes, multiplying this problem of non-  
61 independence, analogous to the problem of pseudoreplication in experimental studies<sup>26-29</sup>.  
62 Phylogenetically-controlled bioinformatic analyses are required to address this problem of  
63 non-independence, and test the robustness of previous conclusions.

64

65 Second, from a theoretical perspective, while horizontal gene transfer can favour the initial  
66 invasion of cooperation, it is not clear if it favours the maintenance of cooperation in the long  
67 run<sup>16</sup>. For example, after a plasmid carrying a cooperative gene has spread through a  
68 population, a loss of function mutation could easily lead to a cheat plasmid evolving, which  
69 could then potentially outcompete the plasmid carrying the cooperative gene<sup>16,30</sup>. Theory is  
70 required that examines the maintenance as well as the invasion of cooperation, while  
71 accounting for important biological details, such as how plasmid transmission depends on the  
72 population frequency of the plasmid, and how frequently plasmids are lost, for example by  
73 segregation during cell division.

74

75 Third, there are alternative hypotheses for why genes coding for extracellular proteins might  
76 be preferentially carried on plasmids in some species (Figure 1)<sup>20,31</sup>. Bacteria can rapidly adapt  
77 to new and/or changing environments by acquiring new genes via horizontal gene transfer, and  
78 losing genes no longer required but costly to maintain (Figure 1b)<sup>32-34</sup>. Genes which facilitate  
79 adaptation to environmental variability are often those which code for molecules secreted  
80 outside the cell<sup>34-37</sup>. Consequently, we might expect to find genes for extracellular proteins on  
81 plasmids to facilitate rapid gain and loss of genes depending on environmental conditions, and  
82 not because they are cooperative *per se*. Alternatively, genes may be favoured to be on plasmids  
83 for reasons other than horizontal gene transfer (Figure 1c)<sup>38</sup>. For example, a higher plasmid  
84 copy number offers a mechanism for more expression of a gene, potentially even conditionally,  
85 in response to certain environmental conditions<sup>38</sup>. The benefit of being able to regulate gene  
86 expression in this way could be higher in genes which code for molecules that are secreted  
87 outside the cell, when different quantities of molecule are required in different environments.  
88 These different hypotheses are not mutually exclusive.

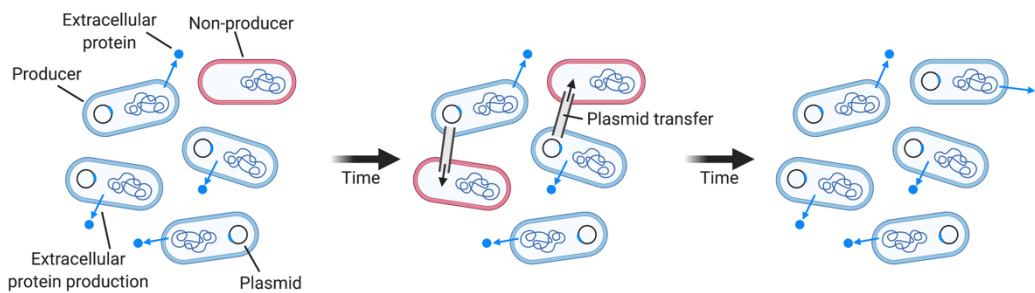
89

90 We addressed all three of these potential problems for the hypothesis that horizontal gene  
91 transfer favours cooperation. We first tested two predictions that would be expected to hold if  
92 horizontal gene transfer favours cooperation. Specifically, cooperative genes would be more  
93 likely to be found on: (i) plasmids relative to chromosomes; (ii) more mobile plasmids relative  
94 to less mobile plasmids<sup>14-20</sup>. We used phylogeny-based statistical methods that control for the  
95 problem of non-independence, analysing 1632 genomes from 51 bacterial species, to examine  
96 the location of genes that code for extracellular proteins. We then used theoretical models, to  
97 examine whether horizontal gene transfer facilitates the evolution as well as the initial spread  
98 of cooperation.

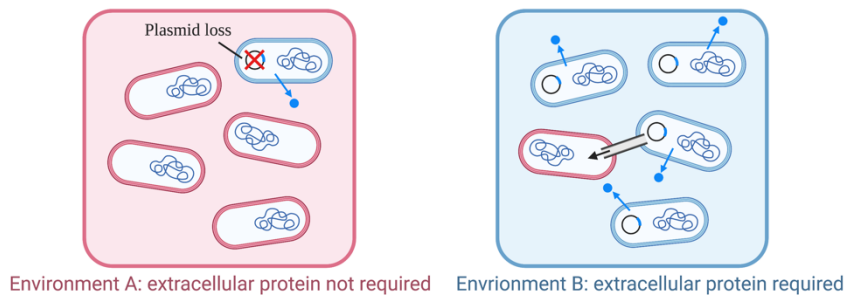
99

100 Finally, we also tested alternative hypotheses for why genes coding for extracellular proteins  
101 might be preferentially carried on plasmids. We used three measures of environmental  
102 variability to ask whether species which had more variable environments were those most  
103 likely to carry genes for extracellular proteins on their plasmids. Additionally, we examined  
104 one of these measures in more detail, to help determine whether genes for extracellular proteins  
105 were located on plasmids so that they could be gained and lost easily (Figure 1b), or instead  
106 because of some additional benefit conferred by plasmid carriage (Figure 1c).  
107

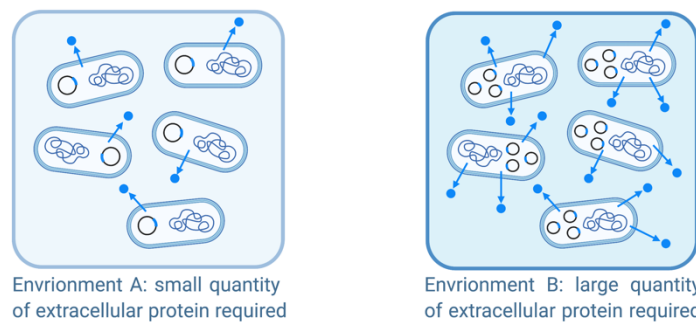
(a) Cooperation Hypothesis: Plasmid transfer stabilises cooperation by 'infecting' non-producing cheats



(b) Gain and Loss Hypothesis: Plasmid transfer allows gain and loss of genes only useful in certain environments



(c) Beyond Horizontal Gene Transfer Hypothesis: Location on plasmid confers advantages beyond mobility



108 **Figure 1. Three hypotheses for why selection might favour genes coding for extracellular**  
109 **proteins to be located on plasmids.**

110 (a) Cooperation Hypothesis. Blue cells produce extracellular proteins which act as cooperative  
111 public goods, while red cells are ‘cheats’ which exploit this cooperation. Over time cheats grow  
112 faster than cooperators since they forgo the cost of public good production. However, because

113 the gene for the extracellular protein is located on a plasmid, cooperators can transfer the gene  
114 to the cheats, turning them into cooperators, increasing genetic relatedness at the cooperative  
115 locus, and stabilising cooperation<sup>14-18</sup>. (b) Gain and Loss Hypothesis. The production of the  
116 extracellular protein is required in some environments, but not others. Transitions between  
117 these environments can result from temporal or spatial change. Cells are selected to either lose  
118 (Environment A) or gain (Environment B) the plasmid coding for the production of the  
119 extracellular protein. (c) Beyond Horizontal Gene Transfer Hypothesis. The location of a gene  
120 on a plasmid could provide a number of benefits, other than the possibility for horizontal gene  
121 transfer<sup>38</sup>. For example, when the quantity of extracellular protein required varies across  
122 environments (A versus B), plasmid copy number could be varied to adjust production<sup>38</sup>.

123

## 124 **Results**

### 125 **Genomic Analyses.**

126 We use the approach developed by Nogueira *et al.*<sup>15,19,20</sup>, of using PSORTb<sup>39</sup> to predict the  
127 subcellular location of every protein encoded by 1632 complete genomes from 51 diverse  
128 bacterial species (Figure S1; Table S3). We are also building upon the work of researchers who  
129 pointed out that extracellular (secreted) proteins are likely to provide a benefit to the local  
130 population of cells, and hence act as cooperative public goods<sup>2,15,19,20,40</sup>. The advantage of this  
131 method is that it allows a large number of genes to be examined, across multiple species.

132

133 Overall, we found the average bacterial genome had 2696 protein-coding genes on the  
134 chromosome(s), and 223 on the plasmid(s). Of these, an average of 57 genes (~2%) coded for  
135 the production of an extracellular protein, with 52 on the chromosome(s) and 5 on the  
136 plasmid(s). This means, on average, 1.9% of chromosome genes and 2.4% of plasmid genes  
137 coded for extracellular proteins. To control for the number of genomes per species, we first  
138 calculated the mean number of genes for each species, and then the mean of these species  
139 means. Therefore, the values above give an indication of the location of genes coding for  
140 extracellular proteins in an average genome. Genes with unknown protein localisations were  
141 not included (Chromosome: 26.2%; Plasmid: 38.3%). Across species, the proportion of genes  
142 coding for extracellular proteins for plasmid(s) was generally more variable than for the  
143 chromosome(s) (Figure S3). These patterns are very similar to those found previously<sup>3,15,19,20</sup>.

144

145 **Extracellular proteins are not overrepresented on plasmids.**

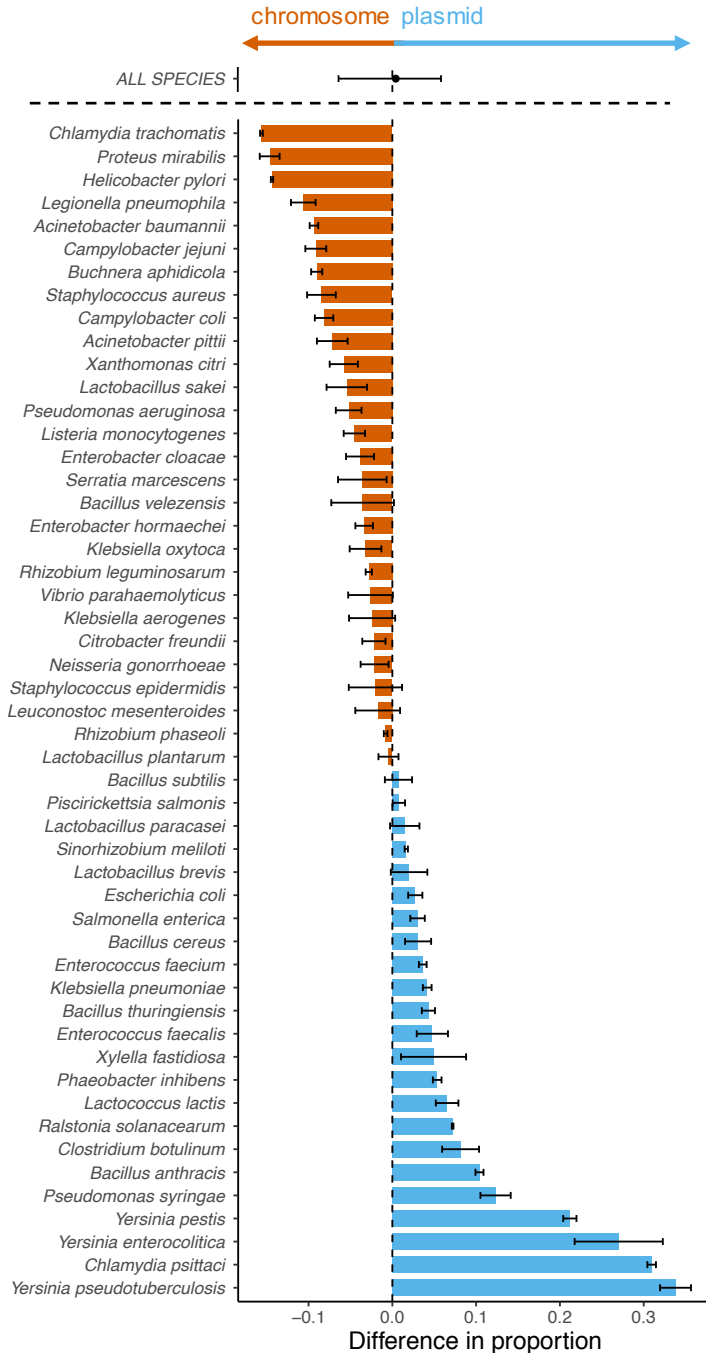
146 We found that extracellular proteins were not more likely to be carried on plasmids compared  
147 to chromosomes (Figure 2). The difference in the proportion of genes that coded for  
148 extracellular proteins between plasmid and chromosome was not significantly different from  
149 zero across all species (MCMCglmm<sup>41</sup>; posterior mean = 0.004, 95% CI = -0.063 to 0.057,  
150 pMCMC= 0.87; n = 1632 genomes; R<sup>2</sup> of species sample size = 0.47, R<sup>2</sup> of phylogeny = 0.17;  
151 Table S2, row 1a). This result was robust to alternative forms of analysis. We also found no  
152 significant difference when we: (i) compared chromosomes to plasmids of only certain  
153 mobilities (Fig S4; Table S2, rows 20-22); (ii) analysed our data by two alternative methods,  
154 by looking at the ratio of proportions instead of the difference, or by considering only whether  
155 the plasmid proportion was greater than the chromosome proportion, removing any effect of  
156 the magnitude of this difference (Figure S5; Table S2, rows 2 and 3). Our analyses use a  
157 bacterial phylogeny, which assumes plasmid evolution follows bacterial phylogeny, but we  
158 also found no significant pattern if we ignored phylogeny and analysed species as independent  
159 data points (Figure 2; Table S2, row 1b; pMCMC = 0.644).

160

161 The lack of an overall significant result was clear when looking at the raw data for the different  
162 species that we examined (Figure 2; Figure S5). There was considerable variation across  
163 species in the location of genes coding for extracellular proteins. Overall, extracellular proteins  
164 were more likely to be on plasmids in 51% of species (26/51), and more likely to be on the  
165 chromosome(s) in 49% (25/51) of species (Figure S5). For example, in *Bacillus anthracis*  
166 genes coding for extracellular proteins were three times more likely to be on plasmids, whereas  
167 in *Acinetobacter baumannii* genes coding for extracellular proteins were three times more  
168 likely to be on the chromosome(s) (Figure S5). Clearly, across species, genes coding for  
169 extracellular proteins are not consistently more likely to be on plasmids.

170

171 As a control, we also analysed the genomic location of the genes coding for all other classes of  
172 protein (Figure S1). Specifically, we analysed genes that coded for the production of  
173 Cytoplasmic, Cytoplasmic Membrane, Periplasmic, Outer Membrane and Cell Wall proteins.  
174 We found that none of these protein localisations were significantly overrepresented on  
175 plasmids or chromosomes across the 51 species (Figure S6; Table S2, rows 5-10). Plasmids  
176 are highly variable in the genes they carry.



177

178

179

180

181

182

183

184

185

**Fig 2. Extracellular proteins are not overrepresented on plasmids.** For each species we calculated the mean difference between plasmid(s) and chromosomes in the proportion of genes coding for extracellular proteins. Species in blue have a difference greater than zero, meaning their plasmid genes code for a greater proportion of extracellular proteins than chromosome genes. Species in red have a difference less than zero, meaning their chromosome genes code for a greater proportion of extracellular proteins than plasmid genes. Error bars indicate the standard error. The dot and error bar at the top of the graph indicate the mean difference and 95% Credible Interval given by a MCMCglmm analysis across all species,

186 controlling for phylogeny and sample size. We arcsine square root transformed proportion data  
187 before calculating the difference. Overall, there is no consistent trend that genes coding for  
188 extracellular proteins are more likely to be carried on plasmids (i.e. no consistent trend towards  
189 species in blue).

190

191 **Importance of controlling for non-independence of genomes.** Our results contrast with  
192 previous studies, which found that plasmid genes code for proportionally more extracellular  
193 proteins than chromosomes<sup>15,19,20</sup>. The first of these studies found this pattern across 20  
194 *Escherichia coli* genomes<sup>15</sup>. We also found that genes coding for extracellular proteins in *E.*  
195 *coli* were more likely to be found on plasmids (Figure 2; Figure S5). However, Figure 2 shows  
196 that this is not a consistent pattern across species: approximately half (25/51) of the species we  
197 analysed showed a pattern in the opposite direction, with genes coding for extracellular proteins  
198 more likely to be on their chromosome(s) than their plasmid(s).

199

200 Two subsequent, multi-species studies found that plasmid genes were significantly more likely  
201 to code for extracellular proteins than chromosome genes<sup>19,20</sup>. These studies used statistical  
202 tests such as Wilcoxon signed-rank test to ask whether there was a consistent pattern, using  
203 bacterial genomes as independent data points. When we analysed our data with the same  
204 statistical methods used in these studies, we also obtained a significant result (Wilcoxon  
205 signed-rank test;  $V = 826530$ ,  $p\text{-value} < 0.001$ ,  $R^2 = 0.385$ ;  $n = 1632$  plasmid-chromosome  
206 pairs). When analysing other questions, Garcia-Garcera & Rocha<sup>20</sup> used MCMCglmm to  
207 control for phylogeny.

208

209 Why does using bacterial genomes as independent data points lead to a significant result? By  
210 using a Wilcoxon signed-rank test, at the level of the genome, we are implicitly assuming that  
211 all the genomes analysed are: (i) independent from one another; (ii) a representative sample of  
212 bacteria in nature. Neither of these are true for multi-species genomic datasets. First, due to  
213 shared ancestry, species are not independent from one another, and so neither are genomes in  
214 such analyses<sup>24,42</sup>. Even a slight lack of independence can lead to heavily biased results in  
215 statistical analyses and spurious conclusions<sup>21</sup>. Second, genomic databases tend to have a  
216 disproportionate abundance of certain species and genera. This will bias the results towards  
217 commonly sequenced species.

218

219 Consequently, when asking questions across species, it is inappropriate to treat all the genomes  
220 in genomic datasets as independent data points. When we performed an analysis analogous to  
221 the Wilcoxon signed-rank test, using the same untransformed data which produced a significant  
222 result above, but controlled for the number of genomes per species and the non-independence  
223 of species, we no longer found any significant difference between the proportion of plasmid  
224 and chromosome genes coding for extracellular proteins (MCMCglimm; posterior mean =  
225 0.017, 95% CI = -0.021 to 0.057, pMCMC = 0.332; n = 1632 plasmid-chromosome paired  
226 differences in extracellular proportion; R<sup>2</sup>: species sample size = 0.46, phylogeny = 0.34; Table  
227 S2, row 4). Furthermore, we found that the number of genomes per species and the non-  
228 independence of species explained 46% and 34% of the variation in data respectively (paired  
229 plasmid and chromosome differences across our 1632 genomes). Taken together, this  
230 illustrates that it is not our data which disagrees with previous studies, but instead our use of  
231 statistical analyses appropriate for multi-genome, multi-species datasets<sup>23-25</sup>.

232

233 These data also illustrate the importance of examining effect sizes, and not just whether results  
234 are statistically significant. With large sample sizes it is possible to get results that are  
235 significant but not biologically important. The percentage of variance explained that is  
236 considered biologically significant can depend upon the kind of data you are examining and  
237 the field of research, but a baseline of 5-10% seems reasonable for many areas of evolutionary  
238 biology (Supp. Info. 1)<sup>43-45</sup>. When bacterial genomes are assumed to be independent data points  
239 in across species analyses, this leads to inflated sample sizes. Consequently, even when results  
240 are statistically significant at P<0.05, they can still only explain 1-2% of the variation in the  
241 data, which is clearly not biologically significant. The flip side of such considerations is that  
242 effects sizes and examination of raw data at the species level (e.g. Figure 2) are also useful  
243 checks against non-significant results due to a lack of statistical power (type II errors).

244

### 245 **Plasmids with higher mobility do not carry more genes for extracellular** 246 **proteins.**

247 We then tested another prediction of the cooperation hypothesis: cooperation is more likely to  
248 be favoured when coded for on more mobile plasmids<sup>14-18</sup>. We used data from the MOBsuite  
249 database to assign plasmids to one of three levels of mobility (Fig 3a)<sup>46,47</sup>. We classify:  
250 conjugative plasmids, which carry all genes necessary to transfer, as the most mobile;  
251 mobilizable plasmids, which are dependent upon conjugative plasmids' machinery to transfer,

252 to have intermediate mobility; non-mobilizable plasmids, which cannot be transferred via  
253 conjugation, to be the least mobile (Fig 3a)<sup>46,48</sup>.

254

255 Genes coding for extracellular proteins were not more likely to be on plasmids with higher  
256 transfer rates (Figure 3b). Examining the slope of the regression between plasmid mobility and  
257 the proportion of genes coding for extracellular proteins, we found no consistent pattern across  
258 species (MCMCglmm; posterior mean = 0.006, 95% CI = -0.040 to 0.052, pMCMC = 0.73; n  
259 = 40; Table S2, row 11). This lack of a significant relationship was robust to different forms of  
260 analysis, including an examination of the means of each mobility type of each species (Figure  
261 S7; Table S2, row 12).

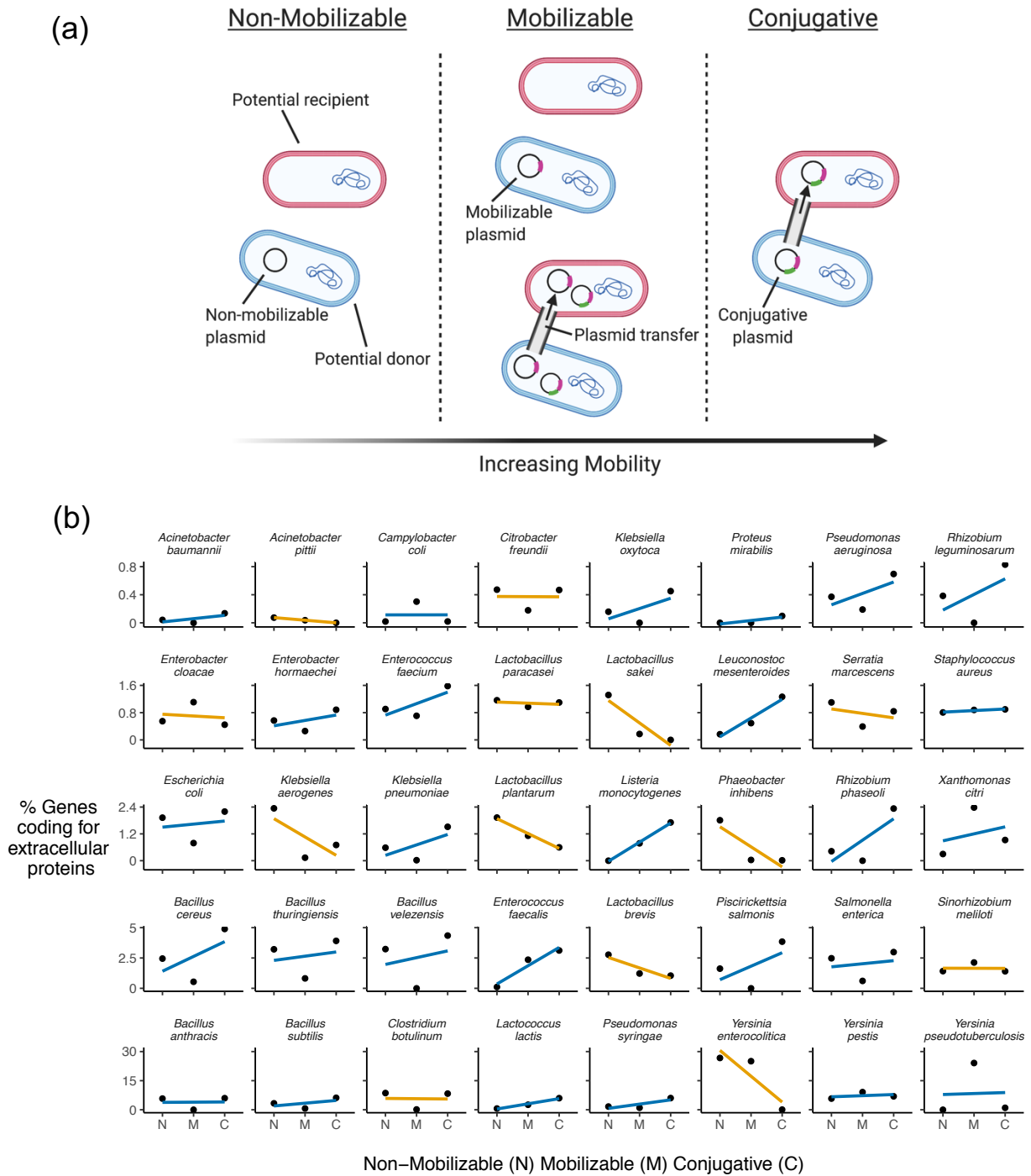
262

263 To examine our assumption that mobilizable plasmids are likely to be less mobile than  
264 conjugative plasmids, we examined how frequently these two kinds of plasmids co-occurred  
265 within a genome. If mobilizable plasmids are present in the same cell as conjugative plasmids,  
266 they could be transmitted at similar rates. However, we found that of genomes with a  
267 mobilizable plasmid(s), 60% did not also carry a conjugative plasmid (434/727). In addition,  
268 when mobilizable plasmids did co-occur with a conjugative plasmid, they did not have a higher  
269 proportion of genes coding for extracellular proteins. (Supp. Info. 1; Figure S10). A caveat here  
270 is that our estimates of transfer rates across different types of plasmid is relative, and it would  
271 be very useful to obtain quantitative estimates of transfer rates.

272

273

274



275 **Figure 3. Plasmid mobility and extracellular proteins.** (a) We divided plasmids into three  
 276 mobility types: non-mobilizable (lowest or no mobility); mobilizable (intermediate mobility);  
 277 conjugative (highest mobility). Blue cells are potential plasmid donors, while red cells are  
 278 potential recipients. Each panel shows when plasmid transfer is possible for one of the three  
 279 plasmid mobility types. Non-mobilizable plasmids cannot be transferred. Mobilizable plasmids  
 280 cannot be transferred alone, but they carry enough genes to ‘hijack’ the machinery of a  
 281 conjugative plasmid that is in the same cell. Conjugative plasmids carry all genes necessary to

282 transfer independently. (b) The 40 species which carried plasmids of all three mobilities are  
283 shown, with a panel for each of these species. Dots in each panel indicate the mean % of genes  
284 coding for extracellular proteins of all plasmids of each mobility level. The lines are the linear  
285 regression of these three points, coloured blue if the slope is positive and orange if the slope is  
286 negative. Note that each row of species has a different y-axis scale, indicated on the left, which  
287 applies to all species in that row. We arcsine square root transformed proportion data before  
288 calculating the mean for each species, and then back-transformed these values for display of  
289 the data. Overall, there is no consistent trend for genes that code for extracellular proteins to  
290 be on more mobile plasmids.

291

## 292 **Theoretical Stability of Cooperation**

293 Our empirical results did not support the theoretical prediction that cooperative genes should  
294 be overrepresented on plasmids, relative to the chromosome<sup>14-18,49</sup>. Consequently, we then  
295 extended existing theory, to examine whether we could find conditions where cooperative  
296 genes were not predicted to be overrepresented on plasmids. We investigated the consequences  
297 of two factors: (1) allowing for a greater range of possible genetic architectures, especially  
298 plasmids that lacked the gene for cooperation (non-cooperative or ‘cheat’ plasmids); and (2)  
299 examining the evolutionary stability (maintenance) of cooperation, not just its initial  
300 invasion<sup>16,49</sup>.

301

302 We examined two possible reasons for why cooperative genes could be overrepresented on  
303 plasmids, relative to the chromosome. First, horizontal gene transfer on a plasmid could allow  
304 cooperation to be favoured in conditions where it would otherwise not be favoured<sup>14-18</sup>. For  
305 example, because plasmid transfer can turn non-cooperators into cooperators, and increase  
306 relatedness at the loci for cooperation<sup>17</sup>. Second, even if horizontal gene transfer did not  
307 increase the range of biological scenarios (parameter space) where cooperation was favoured,  
308 there could be selection for cooperation to be coded for on a plasmid, rather than a  
309 chromosome.

310

311 We assumed an infinite population of haploid individuals (bacterial cells). Individuals may  
312 carry a cooperative gene, that codes for public goods production, either on a plasmid, or the  
313 chromosome, or both (redundancy). We also allowed for the possibility of: non-cooperative  
314 plasmids and chromosomes; plasmid-free cells; a cost of plasmid carriage ( $C_C$ ).

315

316 Each generation, the population is divided into patches, each founded by  $N$  independent cells.  
317 Cells reproduce clonally until there are a large number of cells per patch. Cells are then  
318 randomly shuffled into pairs on their patch and, if a plasmid-free individual has a plasmid-  
319 bearing partner, with probability  $\beta$ , the plasmid-free individual acquires a copy of its partner's  
320 plasmid (horizontal gene transfer). Individuals with a gene for cooperation then produce a  
321 public good, at a cost  $C_G$ , which generates a benefit  $B$  that is shared between all members of  
322 the patch. Individuals then survive according to their fitness. Plasmid-bearing individuals lose  
323 their plasmid with probability  $s$ . Finally, individuals disperse to found new patches.

324

325 Consistent with previous analyses, we found that, in the short term, horizontal gene transfer on  
326 a plasmid can initially help cooperation invade (Figure 4)<sup>14-18</sup>. Horizontal gene transfer  
327 increased the frequency of cooperation, by turning non-cooperators into cooperators, which  
328 also increases relatedness at the cooperative locus on the plasmid<sup>14-18,49</sup>. Relatedness is  
329 increased because, in the short term, whilst plasmids are spreading from rarity, there are many  
330 plasmid-free cells available, meaning plasmids have many opportunities to be transferred,  
331 generating genetic similarity.

332

333 In contrast, we found that transfer on a plasmid did not appreciably increase the range of  
334 parameter space where cooperation was maintained at evolutionary equilibrium (Fig 4a & 5)  
335 (Supp. Info. 4). First, in the absence of plasmid loss ( $s=0$ ), cooperation was only favoured when  
336  $RB-C_G>0$ , where  $R$  is the genetic relatedness at the chromosomal (individual) level ( $R=1/N$ ).  
337 Cooperation was therefore only favoured on the plasmid when it provided a kin selected benefit  
338 at the level of the chromosome (individual), as predicted by Hamilton's rule<sup>50,51</sup>.

339

340 The reason for this result is that, in the absence of plasmid loss ( $s=0$ ), plasmids continue to  
341 increase in frequency after invasion, ultimately reaching fixation in the population. This means  
342 that, in the long term, there are no plasmid-free individuals left to infect, which means that the  
343 overall level of horizontal gene transfer in the population goes to zero. Consequently,  
344 competition between plasmids with and without a cooperative gene (cooperators and cheats)  
345 becomes analogous to the scenario in which the gene for cooperation is on the chromosome<sup>17</sup>.

346

347 Second, when plasmids can be lost ( $s>0$ ), this can favour cooperation on plasmids, but only in  
348 certain areas of parameter space (Figure 5). Plasmid loss means that plasmids do not reach

349 fixation in the population, and so some plasmid transfer still occurs in the evolutionary long  
350 term, increasing relatedness at the cooperative plasmid locus. This increased relatedness may  
351 favour cooperation on the plasmid, when it would not otherwise be favoured on the  
352 chromosome, if plasmids are transferred rapidly (high  $\beta$ ) and rates of plasmid loss are  
353 intermediate (Figure). Specifically, plasmids need to be lost quickly enough that plasmid  
354 relatedness appreciably deviates from chromosomal relatedness, but not too quickly that  
355 plasmids are not maintained (Figure 5). Another factor that might prevent plasmids from  
356 reaching fixation is if there was a constant, high influx of plasmid-free cells (immigration).

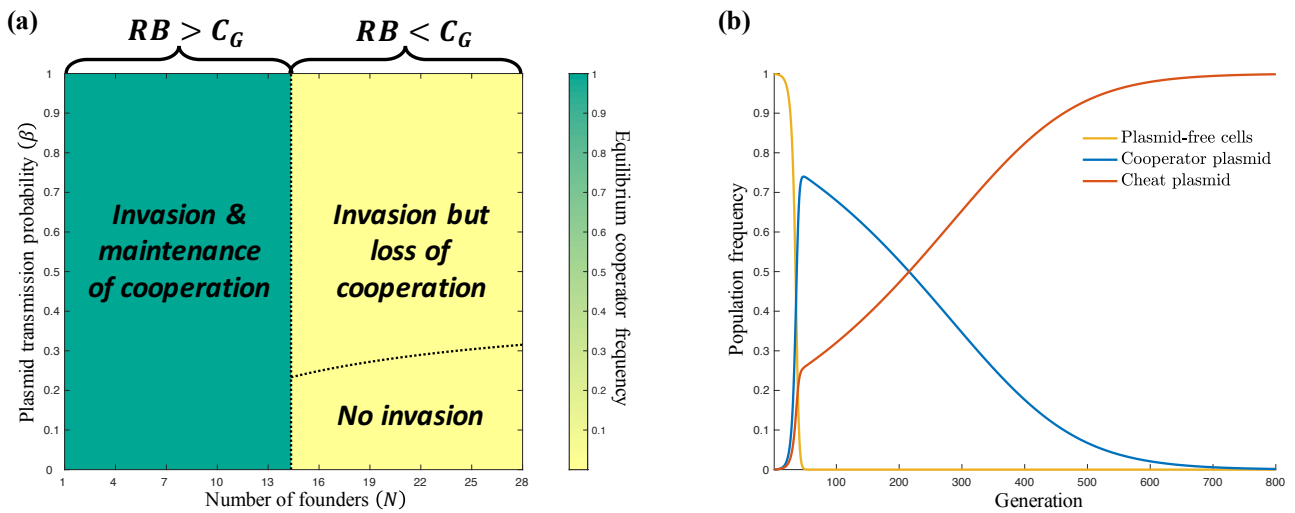
357

358 Overall, our model suggests that horizontal gene transfer can help cooperation initially invade,  
359 but will then often have less influence on whether cooperation is maintained in the long term  
360 (Figures 4 & 5). We are not saying that horizontal gene transfer can never favour cooperation,  
361 just that there is an appreciable area of parameter space where it does not. Consequently, our  
362 model provides an explanation for why cooperative genes are not consistently overrepresented  
363 on plasmids (Figures 2 & 3). An analogous theoretical result for the case without plasmid loss  
364 ( $s=0$ ) was also found in a meta-population model by Mc Ginty *et al.*<sup>16</sup>. Our predictions are  
365 consistent with experiments carried out by Bakkeren *et al.*<sup>30</sup>, who found that location on a  
366 conjugative plasmid could help a cooperative trait invade in *Salmonella* Typhimurium (*S.Tm*),  
367 but that this was only stable with strong population bottlenecks (high relatedness). Dimitriu *et*  
368 *al.*<sup>18</sup> found that cooperative plasmids were favoured in structured but not well-mixed  
369 populations, and that cooperation was favoured more during ‘epidemic spreads’ into a  
370 population.

371

372 In addition, we found that, when cooperation is favoured, cooperative traits are not more likely  
373 to be favoured on, or transferred to, plasmids. The reason is that, when cooperation is favoured,  
374 non-cooperators (cheats) are purged from the population, which means there is no extra fitness  
375 benefit of coding for the cooperative trait on a plasmid rather than the chromosome.  
376 Consequently, our results suggest that horizontal gene transfer only favours cooperation in a  
377 restricted area of parameter space. Although, there could be interesting transient dynamics,  
378 with cooperation being favoured temporarily (Figure 4), or when cooperation has other  
379 consequences, such as increasing plasmid transmission<sup>52,53</sup>. Another important factor is the rate  
380 of horizontal gene transfer. While plasmids clearly transmit fast enough to influence evolution,  
381 the transfer rates per cell per generation might not be high enough to significantly influence  
382 relatedness at the locus for cooperation (i.e. a high enough  $\beta$ )<sup>54</sup>.

383



384

385 **Figure 4. Plasmids facilitate the invasion but not the maintenance of cooperation.** In parts

386 (a) and (b), we plot the results of our theoretical model for the case when there is no plasmid

387 loss ( $s=0$ ). (a) Cooperation is only maintained at equilibrium (green shaded area) when it is

388 favoured at the chromosomal level  $RB > C_G$ , which is unaffected by plasmid transfer ( $\beta$ ). (b)

389 Plasmids can facilitate the invasion and initial spread of cooperation (blue line shoots above

390 red line), but cooperative plasmids are eventually outcompeted by cheat plasmids (red line goes

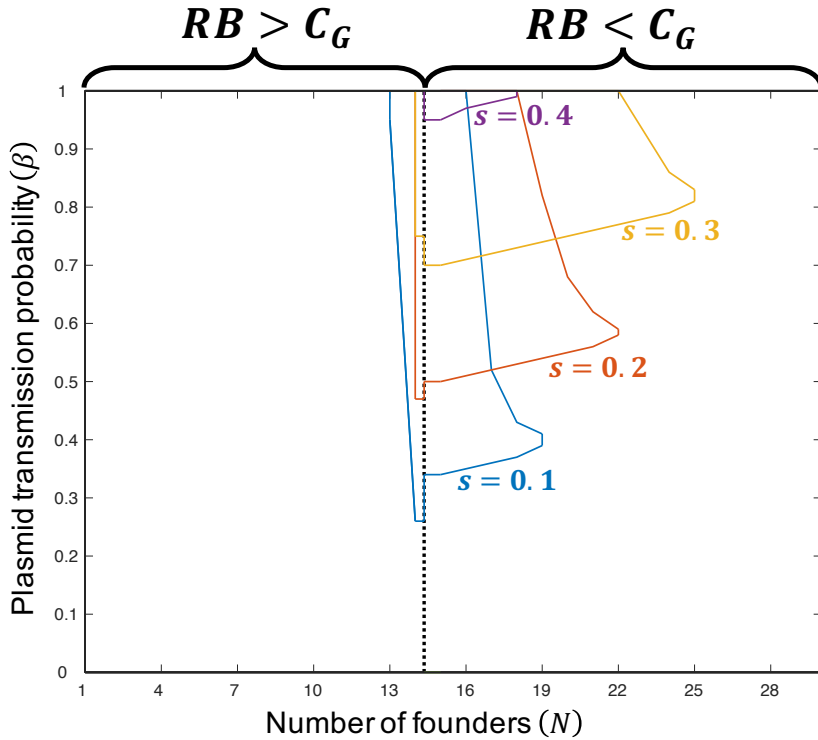
391 to 1). We note that, in (b), all individuals are chromosomal defectors – chromosomal

392 cooperation was permitted, but did not evolve in this run. To generate the plots in (a) and (b),

393 we assumed the following parameter values: (a & b)  $B = 1.435, C_G = 0.1, C_C = 0.2$ ; (b)  $\beta =$

394 0.5,  $N = 16$ .

395



396

397 **Figure 5. Plasmid loss can favour the maintenance of cooperation.** We plot the results of  
 398 our theoretical model for different levels of plasmid loss ( $s=0-1$ ). The areas encapsulated by  
 399 the coloured lines show the regions of parameter space where cooperation is polymorphic at  
 400 equilibrium (i.e. population comprises some cooperators & some defectors). When plasmid  
 401 loss is absent ( $s=0$ ), there is no polymorphism (encapsulated area collapses to nothing),  
 402 meaning cooperation is only maintained at equilibrium (at fixation) when it is favoured at the  
 403 chromosomal level  $RB > C_G$  (to the left of the black dotted line) ( $R=1/N$ ). When plasmid loss  
 404 is intermediate ( $s=0.1,0.2,0.3,0.4$ ), cooperation can be polymorphic at equilibrium  
 405 (encapsulated areas), with cooperation being disfavoured in the encapsulated areas to the left  
 406 of the black dotted line, and favoured in the encapsulated areas to the right of the black dotted  
 407 line, relative to when plasmids are absent ( $\beta=0$ ). When plasmid loss is high ( $s \geq 0.5$ ), or when  
 408 transmission ( $\beta$ ) is low, plasmids fail to persist at equilibrium, meaning they have no long-term  
 409 effect on cooperation (encapsulated areas collapse to nothing). Overall, plasmid loss can  
 410 facilitate cooperation, but only if plasmid loss ( $s$ ) is intermediate and transmission ( $\beta$ ) is high.  
 411 To generate this plot, we assumed the following parameter values:  $B = 1.435, C_G = 0.1, C_C =$   
 412  $0.2$  (same as Fig. 4).

413

414 **Alternate hypotheses**

415 Finally, we examined whether alternate hypotheses may better explain the considerable  
416 variation in the location of genes coding for extracellular proteins across species. Species which  
417 live in more variable environments may be more likely to carry extracellular genes on plasmids.  
418 This could be expected for different reasons, including plasmid transfer allowing genes for  
419 different environments to be gained and lost (Figure 1b), or plasmids conferring some other  
420 advantage not associated with horizontal gene transfer, such as allowing copy number to be  
421 conditionally adjusted (Figure 1c)<sup>31,32,38,55</sup>. There are a number of different ways to classify  
422 environmental variability, and so we used three different methods.

423

424 **Broad host-range pathogens are most likely to carry genes for extracellular proteins on**  
425 **plasmids.** We first used the diversity of pathogen hosts as a proxy for environmental  
426 variability. Although this does not capture all environmental variability experienced by species  
427 in our data set, pathogenicity is a key aspect of bacterial lifestyle that has been suggested to be  
428 important for plasmid gene content, such as antibiotic resistance and virulence factors<sup>6,40,56,57</sup>.  
429 We divided species into three categories: pathogens with broad host-range, pathogens with  
430 narrow host-range, and non-pathogens. Broad host-range pathogens are expected to encounter  
431 more variable environments than narrow host-range pathogens.

432

433 We found that pathogens with a broad host-range were more likely to carry genes coding for  
434 extracellular proteins on their plasmids, compared with both narrow host-range pathogens and  
435 non-pathogens (Fig 6a). Specifically, we compared the difference in the proportion of genes  
436 coding for extracellular proteins between plasmid(s) and chromosome(s) across these three  
437 categories of species (MCMCglmm; Narrow compared to Broad host-range pathogens:  
438 posterior mean = -0.222, 95% CI = -0.322 to -0.123, pMCMC = <0.001; Non-pathogens  
439 compared to Broad host-range pathogens: posterior mean = -0.161, 95% CI = -0.252 to -0.067,  
440 pMCMC = <0.001; n = 701 genomes; R<sup>2</sup> of pathogenicity/host-range = 0.35, R<sup>2</sup> of species  
441 sample size = 0.28, R<sup>2</sup> of phylogeny = 0.11; Table S2, row 23). There was no significant  
442 difference between narrow host-range pathogens and non-pathogens in the proportion of genes  
443 coding for extracellular proteins on their plasmids compared to chromosome(s) (MCMCglmm;  
444 Non-pathogens compared to Narrow host-range pathogens: posterior mean = 0.031, 95% CI =  
445 -0.065 to 0.127, pMCMC = 0.482; n = 389; Table S2, row 25). These patterns hold irrespective  
446 of whether we included species that we could not reliably classify into either category, such as  
447 opportunistic pathogens, in our analyses (Figure S11).

448

449 **Plasmids of broad host-range pathogens carry many pathogenicity genes.** We suspected  
450 that the additional extracellular proteins coded for by plasmids of broad host-range species,  
451 compared to narrow host-range species, may be particularly involved in facilitating  
452 pathogenicity<sup>40,56,57</sup>. To investigate this, we used the program MP3<sup>58</sup> to assign each  
453 extracellular protein as either ‘pathogenic’ or ‘non-pathogenic’.

454

455 We found that plasmids of broad host-range pathogens were particularly enriched with  
456 extracellular proteins involved in facilitating pathogenicity, compared to plasmids of narrow  
457 host-range species (Figure 6b(i)). Specifically, we found that pathogens with a broad host-  
458 range were significantly more likely to code for pathogenic extracellular proteins on their  
459 plasmids compared to narrow host-range species (Figure 6b(i)) (MCMCglmm; Narrow  
460 compared to Broad host-range pathogens: posterior mean = -0.209, 95% CI = -0.350 to -0.086,  
461 pMCMC = 0.012; n=474 genomes; Table S2, row 26). In contrast, the relative location of non-  
462 pathogenic extracellular proteins did not vary between broad and narrow host-range pathogens  
463 (Figure 6b(ii)) (MCMCglmm; Narrow compared to Broad host-range pathogens: posterior  
464 mean = -0.036, 95% CI = -0.115 to 0.040, pMCMC = 0.296; n=474 genomes; Table S2, row  
465 27). Consequently, the excess of genes coding for extracellular proteins on the plasmids of  
466 broad host-range species (Figure 6a) appears to arise due to an excess of pathogenicity genes  
467 coding for extracellular proteins (Figure 6b).

468

469 Most genomic databases are biased towards species that interact with and/or infect humans, so  
470 we examined whether human pathogens had driven the above results. In our dataset, 5 out of  
471 10 broad host-range species and 3 out of 5 narrow host-range species can infect humans. We  
472 found no significant difference in how likely both pathogenic and non-pathogenic extracellular  
473 proteins were to be on plasmids of human pathogens compared to non-human pathogens. We  
474 also found that while host-range had a significant effect on how likely plasmids were to code  
475 for pathogenic extracellular proteins, whether a species could infect humans had no significant  
476 effect (Table S2, rows 28 to 30).

477

478 Pathogenic extracellular proteins could be preferentially coded for on plasmids to facilitate  
479 their gain and loss (Figure 1b: Gain and loss hypothesis), or because of some other benefit  
480 provided by being carried on a plasmid (Figure 1c: Beyond horizontal gene transfer  
481 hypothesis). We tested these possibilities by examining whether pathogenic extracellular  
482 proteins were more likely to be on plasmids that transfer at higher rates. This would be

483 predicted by the gain and loss hypothesis, but not the beyond horizontal gene transfer  
484 hypothesis. We found that plasmids with higher mobility did not code for more pathogenic  
485 extracellular proteins. Specifically, across broad host-range pathogen species, the slope of the  
486 regression between plasmid mobility and the proportion of genes coding for pathogenic  
487 extracellular proteins was not consistently positive (Figure S12) (MCMCglmm; posterior mean  
488 = -0.020, 95% CI = -0.224 to 0.185, pMCMC = 0.774; n=7; Table S2, row 31). This lack of a  
489 significant relationship was robust to additional forms of analysis, such as considering all  
490 pathogenic species, including narrow host-range pathogens and those not carrying plasmids of  
491 all three mobility types (Figure S13; Table S2, rows 32 and 33).

492

493 Taken together, our results are most consistent with the hypothesis that genes coding for  
494 extracellular proteins are overrepresented on plasmids when plasmid carriage provides a  
495 benefit other than mobility (Figure 1c). A number of other factors may influence which genes  
496 are carried on plasmids, beyond horizontal gene transfer. First, there is evidence that increasing  
497 the copy number of plasmids can lead to increasing rates of evolution in the genes they carry<sup>59</sup>,  
498 and it also may act as a mechanism to increase the expression of genes carried on plasmids<sup>60,61</sup>.  
499 For example, increased expression of genes coding for extracellular public goods such as  
500 virulence factors could help invasion of a host and utilisation of host resources. This could be  
501 particularly beneficial for broad host-range pathogens that frequently invade a variety of  
502 different hosts. Copy number of plasmids has also recently been shown to lead to genetic  
503 dominance effects<sup>55</sup>, with likely implications for the phenotypes of genes selected for plasmid  
504 carriage<sup>55</sup>. Second, plasmids compete with their bacterial hosts for resources such as replication  
505 machinery and nucleotides<sup>62,63</sup>. To resolve this competition, plasmids should be under selection  
506 to reduce their cost to the host, with a likely impact on their gene content. For example,  
507 extracellular proteins are, on average, cheaper to produce than intracellular proteins<sup>15,20</sup>.  
508 Plasmid-host competition could consequently select for plasmids to carry more genes coding  
509 for cheaper proteins, and so more extracellular proteins. Our conclusion here should be seen as  
510 tentative, as some form of the gain and loss hypothesis (Figure 1b) could still be argued to be  
511 consistent with the data, if it is just the potential for horizontal gene transfer that matters, and  
512 not the rate.

513



515 **Figure 6. Pathogenicity, host-range and the location of genes coding for extracellular**  
516 **proteins.** We have divided species into either pathogens or non-pathogens, with pathogens  
517 further categorised into those with a narrow or broad host-range. The y-axis in (a) shows the  
518 difference in the proportion of genes on plasmids and chromosomes coding for extracellular  
519 proteins – this is the same as the x-axis in Figure 2. The y-axes in (b)(i) and (b)(ii) show the  
520 difference in the proportion of a subset of genes coding for extracellular proteins on plasmids  
521 and chromosomes which are predicted by MP3 as either (i) pathogenic or (ii) non-pathogenic.  
522 Each dot is the mean for all genomes in a species. Species in blue are those with the relevant  
523 subset of extracellular proteins overrepresented on plasmids, while species in red are those with  
524 the subset of extracellular proteins overrepresented on chromosomes. (c) Phylogeny based on  
525 recently published maximum likelihood tree using 16S ribosomal protein data<sup>64</sup>. The inner ring  
526 indicates whether extracellular proteins were more likely to be coded for on the plasmid(s) or  
527 chromosome(s), as in Figure 2. The outer ring indicates how we classified each species’  
528 pathogenicity, and the presence or absence of diagonal lines for pathogens indicates narrow or  
529 broad host-range, respectively. Species with a pink or green label in the outer ring are those  
530 included in (a) and (b), since for these we could be reasonably confident of whether or not  
531 pathogenicity was an important and consistent aspect of their lifestyle. Overall, pathogens with  
532 a broad host-range are more likely to have genes coding for extracellular proteins, and  
533 particularly those involved in pathogenicity, on their plasmids.

534

535 **Number of environments and core vs accessory genes.** To further examine a potential  
536 association with environmental variability, as could be predicted by both hypotheses b (“Gain  
537 and Loss”) and c (“Beyond Horizontal Gene Transfer”), we also looked at two additional  
538 measures of environmental variability: (i) the number of five broad environments a species was  
539 sequenced in<sup>20,65,66</sup>; (ii) the proportion of a species’ genomes that is composed of ‘core’ genes,  
540 which are those found in all genomes of the species – species which experience more variable  
541 environments appear to have relatively smaller core genomes<sup>32</sup>. We found no significant  
542 correlation between either of these measures and the likelihood that genes coding for  
543 extracellular proteins were carried on plasmids (Figure S14) (Supp. Info. 1; Table S2, rows 35  
544 and 37). Garcia-Garcera & Rocha<sup>20</sup> previously analysed a different but related question,  
545 examining the type of environment, and also used a MCMCglmm to control for the  
546 phylogenetic structure of the data (Supp. Info. 1). Our finding of no correlation between these  
547 two measures of environmental variability and whether plasmids code for extracellular proteins  
548 is in contrast to our above results with respect to pathogen host-range (Figures 5 and 6). This

549 suggests that hypothesis c, which our data is most consistent with, may be important for  
550 pathogens in particular, but not necessarily across all bacterial species and lifestyles.

551

## 552 **Complementary Analyses**

553 There a number of directions in which our analyses could be expanded. We focused on  
554 plasmids because they have been the focus of previous theoretical and empirical work<sup>14,16-18</sup>.  
555 Other mobile genetic elements include bacteriophages and integrative conjugative  
556 elements<sup>67,68</sup>. Comparing core and accessory genes could be a potential way to lump all causes  
557 of horizontal gene transfer<sup>15,19</sup>. We considered the relative transfer rates among mobility types;  
558 quantitative estimates of plasmid transfer rates would be very useful for further examination of  
559 plasmid mobility<sup>48,54,69-71</sup>. We followed previous genomic studies by using extracellular  
560 proteins as indicators of cooperative traits<sup>2,15,19,20</sup>. The advantages of this approach are that: (i)  
561 we could compare our results with those from previous studies; (ii) secretion systems are highly  
562 conserved, allowing us to examine a large number of species, where detailed genetic  
563 annotations are lacking; (iii) cooperation mediated by extracellular proteins is usually  
564 controlled by only one gene, making them potentially more suitable for plasmid carriage  
565 compared to cassettes of multiple genes<sup>72,73</sup>. However, while extracellular proteins are likely  
566 to be cooperative traits, not all cooperative genes code for extracellular proteins (e.g. secondary  
567 metabolites such as siderophores), and not all extracellular proteins are involved in cooperation  
568 (e.g. those involved in motility such as flagellin). It would be very useful to examine more  
569 detailed annotations of social genes, and expand to other mobile genetic elements.

570

## 571 **Discussion**

572 We found no support for the hypothesis that horizontal gene transfer favours cooperation. Our  
573 genomic analyses showed that extracellular proteins are not: (i) overrepresented on plasmids  
574 compared to chromosomes (Figure 2); (ii) more likely to be carried by plasmids that transfer  
575 at higher rates (Figure 3). These patterns could be explained by our theoretical modelling,  
576 which showed that while horizontal gene transfer may help cooperation to initially invade a  
577 population, it has less influence on the maintenance of cooperation in the long term (Figures 4  
578 & 5). Once plasmids become common, cheat plasmids that do not code for cooperation are able  
579 to outcompete cooperative plasmids, analogous to selection at the level of the chromosome<sup>16,30</sup>.  
580 Our results suggest that horizontal gene transfer on plasmids has not consistently favoured  
581 cooperation across bacterial species – but it is still possible that horizontal gene transfer could

582 have an influence in certain scenarios or species. In contrast, we found that genes coding for  
583 extracellular proteins involved in pathogenicity and virulence are preferentially located on  
584 plasmids in pathogens with a broad host-range (Figure 6). These pathogenic virulence genes  
585 were not preferentially located on plasmids that transfer at a higher rate, suggesting that the  
586 benefit of being located on a plasmid is something other than horizontal gene transfer, such as  
587 the ability to vary copy number.

588

## 589 **Methods**

### 590 **Genome Collection**

591 We retrieved 1632 complete genomes comprising 51 bacterial species from GenBank RefSeq  
592 (<https://www.ncbi.nlm.nih.gov>) between February-November 2019. We used species on panX  
593 (<http://pangenome.tuebingen.mpg.de>)<sup>74</sup> as a list of potential species for our dataset, since these  
594 comprise the most sequenced bacterial species. To allow comparison of chromosome and  
595 plasmid genes within the same genome, we only retrieved genomes that contained at least one  
596 plasmid sequence. We included species with 10 or more RefSeq genomes with one or more  
597 plasmids available in our analysis. We retrieved up to 100 genomes for each species; this was  
598 either all complete genomes available for the species, or a random sample where more than  
599 100 were available. Where two or more genomes had the same strain name, we randomly  
600 retrieved one genome to reduce the risk of pseudoreplication.

601

### 602 **Prediction of Subcellular Location of Proteins**

603 We used PSORTb v.3<sup>39</sup> to predict the subcellular location of every protein encoded by each  
604 genome in our dataset. We used a Docker image of PSORTb developed by the Brinkman Lab,  
605 available at: [https://github.com/brinkmanlab/psortb\\_commandline\\_docker](https://github.com/brinkmanlab/psortb_commandline_docker). We chose  
606 PSORTb because it is widely regarded as one of the best performing programs of its kind<sup>75</sup>. It  
607 has also been used in previous analyses to identify ‘cooperative’ genes and/or extracellular  
608 proteins in bacteria<sup>15,20</sup>. The program has a number of modules which are trained to recognise  
609 particular features of proteins. Results from these modules are combined to give a Final  
610 Prediction for each protein. We consulted the literature to confirm the Gram stain of each of  
611 our species. For Gram-positive species, PSORTb assigns proteins to one of four locations  
612 within the cell: cytoplasmic, cytoplasmic membrane, extracellular or cell wall (Figure S1). The  
613 locations for Gram-negative species are the same, except that cell wall is replaced with outer  
614 membrane and periplasmic, meaning there are five possible locations for proteins of Gram-

615 negative species (Figure S1). We used these predicted locations throughout all subsequent  
616 analyses in this work. PSORTb could not reliably assign a subcellular location to 27% of  
617 proteins we analysed, giving a final prediction of ‘unknown’ (Table S1). Unless explicitly  
618 stated, we did not include these unknown proteins in our analyses.

619

## 620 **Predicting Plasmid Mobility**

621 We also predicted the mobility of every plasmid in our dataset using the MOB-typer tool of  
622 the program MOBsuite<sup>46</sup>. This searches for features of plasmid sequences including the origin  
623 of transfer (oriT), relaxase and mating-pair formation to give each plasmid one of three  
624 mobility predictions: (i) conjugative, where plasmids encode all machinery required to transfer  
625 via conjugation; (ii) mobilizable, where plasmids do not encode all machinery, but encode oriT  
626 and/or relaxase, allowing them to ‘hijack’ another plasmid’s conjugation machinery and  
627 mobilize; (iii) non-mobilizable, where plasmids do not encode the genes necessary to be  
628 mobilized by themselves or other plasmids, and so cannot transfer via conjugation. 628 of the  
629 4150 plasmids in our dataset were flagged as ‘unverified’ against the MOBsuite dataset,  
630 meaning their mobility prediction was unreliable and they were not included. This left 3522  
631 plasmids for subsequent analysis.

632

## 633 **Effect of Mobility on Plasmid Extracellular Protein Content**

634 We next examined how plasmid mobility correlates with each plasmid’s extracellular protein  
635 proportion. As part of its mobility prediction, MOBsuite<sup>46</sup> identifies sequences within each  
636 plasmid involved with conjugation. To control for the possibility that conjugative plasmids, by  
637 definition of being conjugative, must carry genes controlling this process, we subtracted the  
638 total number of these sequences from the total number of proteins when calculating the  
639 extracellular proportion of each plasmid. This is a highly conservative control, since it assumes  
640 none of the proteins predicted as extracellular are involved in conjugation. We did all analyses  
641 on these data with and without removing these mating-pair accessions to ensure any results  
642 were not affected by factors unrelated to plasmids’ extracellular protein content.

643

644 Additionally, we used the plasmid mobility predictions to ask whether differences in the  
645 mobility of species’ plasmids correlated with whether genes encoding extracellular proteins  
646 are overrepresented on plasmids compared to chromosomes. We calculated the proportion of  
647 plasmids in each genome capable of transferring via conjugation (conjugative and mobilizable

648 plasmids), and averaged across all genomes to give a general measure of the mobility of each  
649 species' plasmids.

650

### 651 **Measures of Bacterial Lifestyle and Environmental Variability**

652 We classified a species as pathogenic if it was described in the literature as an obligate or  
653 facultative pathogen. Given some bacterial species only rarely act as pathogens, such as  
654 opportunistic pathogens, we only included species where we could be sure pathogenicity was  
655 a key aspect of their lifestyle and a regular selection pressure acting on their genome content.  
656 For this reason, we decided not to include species described as opportunistic pathogens in the  
657 literature and those which frequently live as commensals in their hosts. We classified non-  
658 pathogens as species which are strictly environmental (never live in hosts) or strictly mutualists  
659 and/or commensals (never cause pathogenicity in their hosts). There were 26 species we could  
660 not definitively assign to either of these categories. These were not included in our main  
661 analyses, although we carried out additional analyses to ensure that removing these species did  
662 not bias our results (Figure S10).

663

664 To estimate the host-range of pathogens, we used information from the literature to determine  
665 the maximum taxonomic level of hosts each species is able to invade. We defined narrow host-  
666 range species as those which can invade either only one host species, or host species within the  
667 same genus or family. In contrast, we defined broad-host range pathogens as those capable of  
668 invading host species within the same order, class or phylum. For example, *Xanthomonas citri*  
669 acts as a plant pathogen within the genus *Citrus*<sup>76</sup>, while *Pseudomonas syringae* acts as plant  
670 pathogen across multiple orders of flowering plants<sup>77</sup>. For more details and references to the  
671 literature used for this classification, please see Table S3.

672

673 We completed additional analyses for other two measures and proxies of environmental  
674 variability, the details and results of which can be found in Supp. Info. 1. In brief, we used  
675 previously published data which classified the habitat diversity of species using 16S RNA  
676 environmental datasets across five broad habitats: water, wastewater, sediment, soil and  
677 host<sup>65,66</sup>. We also supplemented this with information from the literature for species not  
678 included in the published data. We used this to ask whether species which lived in multiple  
679 habitats had genes encoding extracellular proteins more overrepresented on their plasmids.

680

681 We also looked at bacterial pangenomes as a proxy for environmental variability, since it has  
682 been noted that species with a high % of accessory genes, defined as genes found in only a  
683 subset of genomes within a species, are generally those with more variable environments. All  
684 pangenome data was collected from panX<sup>74</sup> (<http://pangenome.tuebingen.mpg.de>), since this  
685 calculates the pangenome using the same method across all of our species.

686

### 687 **Pathogenicity categorisation of extracellular proteins**

688 We used MP3<sup>58</sup> to examine the pathogenicity of extracellular protein-coding genes in broad  
689 host-range and narrow host-range pathogens. MP3 compares protein sequences to a curated  
690 dataset of proteins known to be involved in various aspects of pathogenicity: adhesion,  
691 invasion, secretion and resistance<sup>58</sup>. MP3 uses two modules to produce a ‘Hybrid’ prediction  
692 for each protein: either ‘Pathogenic’ or ‘Non-Pathogenic’. We used MP3 with default  
693 parameters to gain this prediction for every extracellular protein in all genomes of broad and  
694 narrow host-range species. MP3 was unable to give a prediction for approximately 9% of  
695 extracellular proteins, and so these were not included in this analysis.

696

697 For each genome in broad and narrow host-range pathogens, we summed the MP3 predictions  
698 to give the total number of ‘Pathogenic’ and ‘Non-Pathogenic’ extracellular proteins on the  
699 chromosome and on the plasmid(s). We then calculated the proportions of plasmid and  
700 chromosome genes which code for ‘Pathogenic’ and ‘Non-Pathogenic’ extracellular proteins.

701

### 702 **Statistical analyses**

703 **MCMCglmm.** Many commonly used statistical methods in biology require data points to be  
704 independent from one another. However, due to shared ancestry, species cannot be considered  
705 as independent data points<sup>24</sup>. Recently developed statistical methods now allow for  
706 phylogenetic relationships to be controlled for within mixed effects models. For all statistical  
707 analyses we used the MCMCglmm (Markov Chain Monte Carlo generalised linear mixed  
708 effects model) package in R with phylogeny a random effect<sup>41,78</sup>. This means the phylogeny is  
709 implemented in the model as a covariance matrix of the relationships between species, which  
710 is controlled for when considering whether patterns exist across species<sup>41,78</sup>. We also included  
711 sample size as a random effect when analysing at the genome level to control for differences  
712 in the number of genomes per species. Specific details of each model can be found in Table  
713 S2. We extracted from each model the posterior mean, 95% Credible Intervals (functionally

714 similar to 95% Confidence Intervals), and the pMCMC value (generally interpreted in a similar  
715 way to a ‘p-value’). We also calculated  $R^2$  values for models of particular interest using  
716 methods described in<sup>79,80</sup>. A detailed description of MCMCglmm can be found elsewhere<sup>41,78</sup>.

717

718 The response variable in all of our analyses is either a proportion or a measure calculated from  
719 proportions. Proportion data is bound between 0 and 1 and has a non-normal distribution. To  
720 control for this, all proportion data in our analyses has been arcsine square root transformed to  
721 improve normality.

722

723 **Phylogeny.** To control for species relationships, we generated a phylogeny including all 51  
724 species in our dataset (Fig S2). We used a recently published maximum likelihood tree using  
725 16S ribosomal protein data as the basis for our phylogeny<sup>64</sup>. This tree of life typically had only  
726 one representative species per genus. We used the R package ‘ape’ to extract all branches  
727 matching species in our dataset<sup>81</sup>. In cases where the genus representative was different to the  
728 species in our dataset, we swapped the tip name with our species, since all members of the  
729 same genus are equally related to members of a sister genus. In cases where we had multiple  
730 species within a single genus in our dataset, we used the R package ‘phylotools’ to add these  
731 species as additional branches into their genus<sup>82</sup>. We used published phylogenies from the  
732 literature to add any within-genus clustering of species’ branches. We used this phylogeny in  
733 nexus format for all our MCMCglmm analyses (Fig S2, Table S2). Methods are also available  
734 to control for uncertainty in phylogenetic reconstruction<sup>83,84</sup>, although we have not done this  
735 here.

736

### 737 **Acknowledgements**

738 We thank: Craig MacLean, Kevin Foster, Laurence Belcher, Chunhui Hao, and especially  
739 Eduardo Rocha for their helpful comments; James Robertson for providing plasmid mobility  
740 data from the MOBSuite database; the BBSRC (A.E.D.), ERC (J.L.T., T.W.S., A.S.G., M.G.  
741 and S.A.W.), and NSERC-CRSNG of Canada (G.W.) for funding. We also thank Alex  
742 Washburne and three anonymous reviewers for comments which greatly improved the  
743 manuscript. Conceptual figures were created with Biorender.com.

744

### 745 **Author Contributions**

746 A.E.D., J.L.T., A.S.G., S.A.W and M.G. conceived the genomic analyses and interpreted  
747 results. A.E.D. and J.L.T. collected and analysed genomic data, and A.E.D. produced the  
748 corresponding figures. T.W.S, G.W. and S.A.W. conceived the theoretical modelling and  
749 interpreted results. T.W.S. completed the formal theoretical modelling. A.E.D., J.L.T, T.W.S.,  
750 S.A.W., and M.G. wrote and/or edited the manuscript. A.E.D. wrote and put together S1, S2  
751 and S3, and T.W.S. wrote and put together S4. All authors commented on and approved the  
752 manuscript for submission.

753

### 754 **Competing Interests**

755 The authors declare no competing interests.

756

### 757 **Data Availability Statement**

758 The datasets generated and/or analysed (including accession codes) during the current study  
759 are available from the corresponding author on request, and will be made available when  
760 published.

761

### 762 **Code Availability Statement**

763 Code for the bioinformatic/statistical analyses and simulations are available from the  
764 corresponding author on request.

765

### 766 **References**

- 767 1. Foster, K. R. Social behaviour in microorganisms. in *Social Behaviour* (eds. Szekely, T.,  
768 Moore, A. J. & Komdeur, J.) 331–356 (Cambridge University Press, 2010).  
769 doi:10.1017/CBO9780511781360.027.
- 770 2. McNally, L., Viana, M. & Brown, S. P. Cooperative secretions facilitate host range  
771 expansion in bacteria. *Nat. Commun.* **5**, (2014).
- 772 3. West, S. A., Griffin, A. S., Gardner, A. & Diggle, S. P. Social evolution theory for  
773 microorganisms. *Nat. Rev. Microbiol.* **4**, 597–607 (2006).
- 774 4. Simonet, C. & McNally, L. Kin selection explains the evolution of cooperation in the gut  
775 microbiota. *Proc. Natl. Acad. Sci.* **118**, (2021).
- 776 5. Griffin, A. S., West, S. A. & Buckling, A. Cooperation and competition in pathogenic  
777 bacteria. *Nature* **430**, 1024–1027 (2004).
- 778 6. Hale, T. L. Genetic basis of virulence in *Shigella* species. *Microbiol. Rev.* **55**, 206–224  
779 (1991).
- 780 7. Dinges, M. M., Orwin, P. M. & Schlievert, P. M. Exotoxins of *Staphylococcus aureus*.  
781 *Clin. Microbiol. Rev.* **13**, 16–34, table of contents (2000).

- 782 8. Diggle, S. P., Griffin, A. S., Campbell, G. S. & West, S. A. Cooperation and conflict in  
783 quorum-sensing bacterial populations. *Nature* **450**, 411–414 (2007).
- 784 9. Jones, S. *et al.* The lux autoinducer regulates the production of exoenzyme virulence  
785 determinants in *Erwinia carotovora* and *Pseudomonas aeruginosa*. *EMBO J.* **12**, 2477–  
786 2482 (1993).
- 787 10. Sandoz, K. M., Mitzimberg, S. M. & Schuster, M. Social cheating in *Pseudomonas*  
788 *aeruginosa* quorum sensing. *Proc. Natl. Acad. Sci.* **104**, 15876–15881 (2007).
- 789 11. Ghoul, M., Griffin, A. S. & West, S. A. Toward an evolutionary definition of cheating.  
790 *Evolution* **68**, 318–331 (2014).
- 791 12. Butaitė, E., Baumgartner, M., Wyder, S. & Kümmerli, R. Siderophore cheating and  
792 cheating resistance shape competition for iron in soil and freshwater *Pseudomonas*  
793 communities. *Nat. Commun.* **8**, 414 (2017).
- 794 13. Thomas, C. & Nielsen, K. Thomas CM, Nielsen KM. Mechanisms of, and barriers to,  
795 horizontal gene transfer between bacteria. *Nat Rev Micro* 3: 711-721. *Nat. Rev.*  
796 *Microbiol.* **3**, 711–21 (2005).
- 797 14. Smith, J. The social evolution of bacterial pathogenesis. *Proc. R. Soc. Lond. B Biol. Sci.*  
798 **268**, 61–69 (2001).
- 799 15. Nogueira, T. *et al.* Horizontal Gene Transfer of the Secretome Drives the Evolution of  
800 Bacterial Cooperation and Virulence. *Curr. Biol.* **19**, 1683–1691 (2009).
- 801 16. Mc Ginty, S. E., Rankin, D. J. & Brown, S. P. Horizontal gene transfer and the evolution  
802 of bacterial cooperation: mobile elements and bacterial cooperation. *Evolution* **65**, 21–32  
803 (2011).
- 804 17. Mc Ginty, S. É., Lehmann, L., Brown, S. P. & Rankin, D. J. The interplay between  
805 relatedness and horizontal gene transfer drives the evolution of plasmid-carried public  
806 goods. *Proc. R. Soc. B Biol. Sci.* **280**, 20130400 (2013).
- 807 18. Dimitriu, T. *et al.* Genetic information transfer promotes cooperation in bacteria. *Proc.*  
808 *Natl. Acad. Sci.* **111**, 11103–11108 (2014).
- 809 19. Nogueira, T., Touchon, M. & Rocha, E. P. C. Rapid Evolution of the Sequences and  
810 Gene Repertoires of Secreted Proteins in Bacteria. *PLoS ONE* **7**, e49403 (2012).
- 811 20. Garcia-Garcera, M. & Rocha, E. P. C. Community diversity and habitat structure shape  
812 the repertoire of extracellular proteins in bacteria. *Nat. Commun.* **11**, 758 (2020).
- 813 21. Kruskal, W. Miracles and Statistics: The Casual Assumption of Independence. *J. Am.*  
814 *Stat. Assoc.* **83**, 929–940 (1988).
- 815 22. Ives, A. R. & Zhu, J. Statistics for correlated data: phylogenies, space, and time. *Ecol.*  
816 *Appl. Publ. Ecol. Soc. Am.* **16**, 20–32 (2006).
- 817 23. Felsenstein, J. Phylogenies and the Comparative Method. *Am. Nat.* **125**, 1–15 (1985).
- 818 24. Harvey, P. H. & Pagel, M. D. *The Comparative Method in Evolutionary Biology.* (Oxford  
819 University Press, 1991).
- 820 25. Grafen, A. The phylogenetic regression. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **326**,  
821 119–157 (1989).
- 822 26. Hurlbert, S. H. Pseudoreplication and the Design of Ecological Field Experiments. *Ecol.*  
823 *Monogr.* **54**, 187–211 (1984).
- 824 27. Ruxton, G. & Colegrave, N. *Experimental Design for the Life Sciences.* (OUP Oxford,  
825 2011).

- 826 28. Stone, G. N., Nee, S. & Felsenstein, J. Controlling for non-independence in comparative  
827 analysis of patterns across populations within species. *Philos. Trans. R. Soc. B Biol. Sci.*  
828 **366**, 1410–1424 (2011).
- 829 29. Ives, A. R., Midford, P. E. & Garland, T., Jr. Within-Species Variation and Measurement  
830 Error in Phylogenetic Comparative Methods. *Syst. Biol.* **56**, 252–270 (2007).
- 831 30. Bakkeren, E. *et al.* Cooperative virulence can emerge via horizontal gene transfer but is  
832 stabilized by transmission. *bioRxiv* 2021.02.11.430745 (2021)  
833 doi:10.1101/2021.02.11.430745.
- 834 31. Ghouil, M., Andersen, S. B. & West, S. A. Sociomics: Using Omic Approaches to  
835 Understand Social Evolution. *Trends Genet.* **33**, 408–419 (2017).
- 836 32. McInerney, J. O., McNally, A. & O’Connell, M. J. Why prokaryotes have pangenomes.  
837 *Nat. Microbiol.* **2**, 17040 (2017).
- 838 33. Niehus, R., Mitri, S., Fletcher, A. G. & Foster, K. R. Migration and horizontal gene  
839 transfer divide microbial genomes into multiple niches. *Nat. Commun.* **6**, (2015).
- 840 34. Cordero, O. X. *et al.* Ecological Populations of Bacteria Act as Socially Cohesive Units  
841 of Antibiotic Production and Resistance. *Science* **337**, 1228–1231 (2012).
- 842 35. Rakoff-Nahoum, S., Coyne, M. J. & Comstock, L. E. An Ecological Network of  
843 Polysaccharide Utilization among Human Intestinal Symbionts. *Curr. Biol.* **24**, 40–49  
844 (2014).
- 845 36. Nocelli, N., Bogino, P. C., Banchio, E. & Giordano, W. Roles of Extracellular  
846 Polysaccharides and Biofilm Formation in Heavy Metal Resistance of Rhizobia.  
847 *Materials* **9**, 418 (2016).
- 848 37. Ciofu, O., Beveridge, T. J., Kadurugamuwa, J., Walther-Rasmussen, J. & Høiby, N.  
849 Chromosomal  $\beta$ -lactamase is packaged into membrane vesicles and secreted from  
850 *Pseudomonas aeruginosa*. *J. Antimicrob. Chemother.* **45**, 9–13 (2000).
- 851 38. Rodríguez-Beltrán, J., DelaFuente, J., León-Sampedro, R., MacLean, R. C. & San Millán,  
852 Á. Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. *Nat. Rev.*  
853 *Microbiol.* 1–13 (2021) doi:10.1038/s41579-020-00497-1.
- 854 39. Yu, N. Y. *et al.* PSORTb 3.0: improved protein subcellular localization prediction with  
855 refined localization subcategories and predictive capabilities for all prokaryotes.  
856 *Bioinformatics* **26**, 1608–1615 (2010).
- 857 40. Rankin, D. J., Rocha, E. P. C. & Brown, S. P. What traits are carried on mobile genetic  
858 elements, and why? *Heredity* **106**, 1–10 (2011).
- 859 41. Hadfield, J. D. MCMC Methods for Multi-Response Generalized Linear Mixed Models:  
860 The MCMCglmm R Package. *J. Stat. Softw.* **33**, 1–22 (2010).
- 861 42. Clutton-Brock, T. H. & Harvey, P. H. Primate ecology and social organization. *J. Zool.*  
862 **183**, 1–39 (1977).
- 863 43. Jennions, M. D. & Møller, A. P. A survey of the statistical power of research in  
864 behavioral ecology and animal behavior. *Behav. Ecol.* **14**, 438–445 (2003).
- 865 44. Crawley, M. J. *Statistics: An Introduction Using R*. (John Wiley & Sons, 2014).
- 866 45. Cohen, J. *Statistical Power Analysis for the Behavioral Sciences*. (Routledge, 1988).
- 867 46. Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering, reconstruction  
868 and typing of plasmids from draft assemblies. *Microb. Genomics* **4**, (2018).

- 869 47. Robertson, J., Bessonov, K., Schonfeld, J. & Nash, J. H. E. Universal whole-sequence-  
870 based plasmid typing and its utility to prediction of host range and epidemiological  
871 surveillance. *Microb. Genomics* **6**, (2020).
- 872 48. Smillie, C., Garcillan-Barcia, M. P., Francia, M. V., Rocha, E. P. C. & de la Cruz, F.  
873 Mobility of Plasmids. *Microbiol. Mol. Biol. Rev.* **74**, 434–452 (2010).
- 874 49. Mc Ginty, S. É. & Rankin, D. J. The evolution of conflict resolution between plasmids  
875 and their bacterial hosts. *Evolution* **66**, 1662–1670 (2012).
- 876 50. Hamilton, W. D. Genetical evolution of social behaviour I & II. *J Theor Biol* **7**, 1–52  
877 (1964).
- 878 51. work(s);, W. D. H. R. The Evolution of Altruistic Behavior. *Am. Nat.* **97**, 354–356  
879 (1963).
- 880 52. Ghigo, J. M. Natural conjugative plasmids induce bacterial biofilm development. *Nature*  
881 **412**, 442–445 (2001).
- 882 53. Di Venanzio, G. *et al.* Multidrug-resistant plasmids repress chromosomally encoded  
883 T6SS to enable their dissemination. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 1378–1383  
884 (2019).
- 885 54. Sheppard, R. J., Beddis, A. E. & Barraclough, T. G. The role of hosts, plasmids and  
886 environment in determining plasmid transfer rates: A meta-analysis. *Plasmid* **108**,  
887 102489 (2020).
- 888 55. Rodríguez-Beltrán, J. *et al.* Genetic dominance governs the evolution and spread of  
889 mobile genetic elements in bacteria. *Proc. Natl. Acad. Sci.* **117**, 15755–15762 (2020).
- 890 56. Cornelis, G. R. *et al.* The Virulence Plasmid of Yersinia, an Antihost Genome. *Microbiol.*  
891 *Mol. Biol. Rev.* **62**, 1315–1352 (1998).
- 892 57. Köstlbacher, S., Collingro, A., Halter, T., Domman, D. & Horn, M. Coevolving Plasmids  
893 Drive Gene Flow and Genome Plasticity in Host-Associated Intracellular Bacteria. *Curr.*  
894 *Biol.* **31**, 346-357.e3 (2021).
- 895 58. Gupta, A., Kapil, R., Dhakan, D. B. & Sharma, V. K. MP3: A Software Tool for the  
896 Prediction of Pathogenic Proteins in Genomic and Metagenomic Data. *PLOS ONE* **9**,  
897 e93907 (2014).
- 898 59. San Millan, A., Escudero, J. A., Gifford, D. R., Mazel, D. & MacLean, R. C. Multicopy  
899 plasmids potentiate the evolution of antibiotic resistance in bacteria. *Nat. Ecol. Evol.* **1**,  
900 0010 (2017).
- 901 60. Carrier, T., Jones, K. L. & Keasling, J. D. mRNA stability and plasmid copy number  
902 effects on gene expression from an inducible promoter system. *Biotechnol. Bioeng.* **59**,  
903 666–672 (1998).
- 904 61. Rodríguez-Beltrán, J. *et al.* Multicopy plasmids allow bacteria to escape from fitness  
905 trade-offs during evolutionary innovation. *Nat. Ecol. Evol.* **2**, 873–881 (2018).
- 906 62. Dietel, A.-K., Kaltenpoth, M. & Kost, C. Convergent Evolution in Intracellular Elements:  
907 Plasmids as Model Endosymbionts. *Trends Microbiol.* **26**, 755–768 (2018).
- 908 63. Rocha, E. P. C. & Danchin, A. Base composition bias might result from competition for  
909 metabolic resources. *Trends Genet.* **18**, 291–294 (2002).
- 910 64. Hug, L. A. *et al.* A new view of the tree of life. *Nat. Microbiol.* **1**, 16048 (2016).

- 911 65. Garcia-Garcera, M., Touchon, M., Brisse, S. & Rocha, E. P. C. Metagenomic assessment  
912 of the interplay between the environment and the genetic diversification of  
913 *Acinetobacter*. *Environ. Microbiol.* **19**, 5010–5024 (2017).
- 914 66. Kümmerli, R., Schiessl, K. T., Waldvogel, T., McNeill, K. & Ackermann, M. Habitat  
915 structure and the evolution of diffusible siderophores in bacteria. *Ecol. Lett.* **17**, 1536–  
916 1544 (2014).
- 917 67. Canchaya, C., Fournous, G., Chibani-Chennoufi, S., Dillmann, M. L. & Brüssow, H.  
918 Phage as agents of lateral gene transfer. *Curr. Opin. Microbiol.* **6**, 417–424 (2003).
- 919 68. Burrus, V. & Waldor, M. K. Shaping bacterial genomes with integrative and conjugative  
920 elements. *Res. Microbiol.* **155**, 376–386 (2004).
- 921 69. O’Brien, F. G. *et al.* Origin-of-transfer sequences facilitate mobilisation of non-  
922 conjugative antimicrobial-resistance plasmids in *Staphylococcus aureus*. *Nucleic Acids*  
923 *Res.* **43**, 7971–7983 (2015).
- 924 70. Rodríguez-Rubio, L. *et al.* Extensive antimicrobial resistance mobilization via multicopy  
925 plasmid encapsidation mediated by temperate phages. *J. Antimicrob. Chemother.* **75**,  
926 3173–3180 (2020).
- 927 71. Ramsay, J. P. & Firth, N. Diverse mobilization strategies facilitate transfer of non-  
928 conjugative mobile genetic elements. *Curr. Opin. Microbiol.* **38**, 1–9 (2017).
- 929 72. Jain, R., Rivera, M. C. & Lake, J. A. Horizontal gene transfer among genomes: the  
930 complexity hypothesis. *Proc. Natl. Acad. Sci. U. S. A.* **96**, 3801–3806 (1999).
- 931 73. Cohen, O., Gophna, U. & Pupko, T. The complexity hypothesis revisited: connectivity  
932 rather than function constitutes a barrier to horizontal gene transfer. *Mol. Biol. Evol.* **28**,  
933 1481–1489 (2011).
- 934 74. Ding, W., Baumdicker, F. & Neher, R. A. panX: pan-genome analysis and exploration.  
935 *Nucleic Acids Res.* **46**, e5 (2018).
- 936 75. Gardy, J. L. & Brinkman, F. S. L. Methods for predicting bacterial protein subcellular  
937 localization. *Nat. Rev. Microbiol.* **4**, 741–751 (2006).
- 938 76. Ference, C. M. *et al.* Recent advances in the understanding of *Xanthomonas citri* ssp. *citri*  
939 pathogenesis and citrus canker disease management. *Mol. Plant Pathol.* **19**, 1302–1318  
940 (2018).
- 941 77. Morris, C. E., Lamichhane, J. R., Nikolić, I., Stanković, S. & Moury, B. The overlapping  
942 continuum of host range among strains in the *Pseudomonas syringae* complex.  
943 *Phytopathol. Res.* **1**, 4 (2019).
- 944 78. Hadfield, J. D. MCMCglmm Course Notes. Available at [cran.us.r-](http://cran.us.r-project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf)  
945 [project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf](http://cran.us.r-project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf). (2019).
- 946 79. Nakagawa, S. & Schielzeth, H. A general and simple method for obtaining R<sup>2</sup> from  
947 generalized linear mixed-effects models. *Methods Ecol. Evol.* **4**, 133–142 (2013).
- 948 80. Nakagawa, S., Johnson, P. C. D. & Schielzeth, H. The coefficient of determination R<sup>2</sup>  
949 and intra-class correlation coefficient from generalized linear mixed-effects models  
950 revisited and expanded. *J R Soc Interface* **11** (2017).
- 951 81. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and  
952 evolutionary analyses in R. *Bioinforma. Oxf. Engl.* **35**, 526–528 (2019).
- 953 82. Revell, L. J. phytools: an R package for phylogenetic comparative biology (and other  
954 things). *Methods Ecol. Evol.* **3**, 217–223 (2012).

- 955 83. Washburne, A. D. *et al.* Methods for phylogenetic analysis of microbiome data. *Nat.*  
956 *Microbiol.* **3**, 652–661 (2018).  
957 84. Som, A. Causes, consequences and solutions of phylogenetic incongruence. *Brief.*  
958 *Bioinform.* **16**, 536–548 (2015).  
959