

1 **Plasmids facilitate pathogenicity, not cooperation, in bacteria**

2 Anna E. Dewar^{1,a,*}, Joshua L. Thomas^{1,a}, Thomas W. Scott¹, Geoff Wild²,
3 Ashleigh S. Griffin¹, Stuart A. West^{1,b}, Melanie Ghoul^{1,b}

4 ¹Department of Zoology, University of Oxford, Oxford, OX1 3SZ, United Kingdom

5 ²Department of Applied Mathematics, University of Western Ontario, London, Ontario N6A
6 3K7, Canada

7 a Joint first author

8 b Joint last author

9 *Corresponding author

10

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

13

14 **Abstract**

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

27

28

29

30

31

32 **Introduction**

33 The growth and success of many bacterial populations depends upon the production of
34 cooperative ‘public goods’¹⁻⁴. Public goods are molecules whose secretion provides a benefit
35 to the local group of cells. Examples include iron-scavenging siderophores⁵, exotoxins that
36 disintegrate host cell membranes^{6,7}, and elastases that break down connective tissues⁸⁻¹⁰. A
37 problem is that cooperation can be exploited by ‘cheats’: cells which avoid the cost of
38 producing public goods but can still use and benefit from those produced by cooperative
39 cells^{3,11,12}. What prevents cheats from outcompeting cooperators, and ultimately destabilising
40 cooperation?

41
42 In bacteria, some genetic elements are able to move between cells¹³. This horizontal gene
43 transfer has been suggested as a mechanism to help stabilize the production of cooperative
44 public goods¹⁴⁻¹⁸ (Figure 1a). If a gene coding for the production of a public good can be
45 transferred horizontally, it would allow cheats to be ‘infected’ with the cooperative gene and
46 turned into cooperators, increasing genetic relatedness at the cooperative locus. Theoretical
47 models have shown that this can facilitate the invasion of cooperative genes, in conditions
48 where they would not be favoured on chromosomes¹⁴⁻¹⁸. Experiments have supported this
49 prediction¹⁸. In addition, bioinformatic analyses across a range of species found that genes that
50 code for extracellular proteins, many of which act as public goods, are more likely to be found
51 on plasmids than the chromosome^{15,19,20}.

52
53 There are, however, three potential problems for the hypothesis that horizontal gene transfer
54 favours cooperation. First, previous bioinformatic analyses made important first steps, but are
55 not conclusive. One study examined only a single species, which may not be representative of
56 all bacteria¹⁵. Two additional studies examined multiple species, but assumed that genes and
57 genomes from the same and different species can be treated as independent data points, in a
58 way that could have led to spurious results^{19,20}. Statistical tests typically assume that data points
59 are independent, and even slight non-independence can lead to heavily biased results (type I
60 errors)^{21,22}. There is an extensive literature in the field of evolutionary biology showing that
61 species share characteristics inherited through common descent, rather than through
62 independent evolution, and so cannot be considered independent data points²³⁻²⁵. Genomes are
63 nested within species, and genes are nested within genomes, multiplying this problem of non-
64 independence, analogous to the problem of pseudoreplication in experimental studies²⁶⁻²⁹.

65 Phylogenetically-controlled bioinformatic analyses are required to address this problem of
66 non-independence, and test the robustness of previous conclusions.

67

68 Second, from a theoretical perspective, while horizontal gene transfer can favour the initial
69 invasion of cooperation, it is not clear if it favours the maintenance of cooperation in the long
70 run¹⁶. For example, after a plasmid carrying a cooperative gene has spread through a
71 population, a loss of function mutation could easily lead to a cheat plasmid evolving, which
72 could then potentially outcompete the plasmid carrying the cooperative gene^{16,30}. Theory is
73 required that examines the maintenance as well as the invasion of cooperation, while
74 accounting for important biological details, such as how plasmid transmission depends on the
75 population frequency of the plasmid.

76

77 Third, there are alternative hypotheses for why genes coding for extracellular proteins might
78 be preferentially carried on plasmids in some species (Figure 1)^{20,31}. Bacteria can rapidly adapt
79 to new and/or changing environments by acquiring new genes via horizontal gene transfer, and
80 losing genes no longer required but costly to maintain (Figure 1b)³²⁻³⁴. Genes which facilitate
81 adaptation to environmental variability are often those which code for molecules secreted
82 outside the cell³⁴⁻³⁷. Consequently, we might expect to find genes for extracellular proteins on
83 plasmids to facilitate rapid gain and loss of genes depending on environmental conditions, and
84 not because they are cooperative *per se*. Alternatively, genes may be favoured to be on plasmids
85 for reasons other than horizontal gene transfer (Figure 1c)³⁸. For example, a higher plasmid
86 copy number offers a mechanism for more expression of a gene, potentially even conditionally,
87 in response to certain environmental conditions³⁸. The benefit of being able to regulate gene
88 expression in this way could be higher in genes which code for molecules that are secreted
89 outside the cell, when different quantities of molecule are required in different environments.

90

91 We addressed all three of these potential problems for the hypothesis that horizontal gene
92 transfer favours cooperation. We first tested two predictions that would be expected to hold if
93 horizontal gene transfer favours cooperation. Specifically, cooperative genes would be more
94 likely to be found on: (i) plasmids relative to chromosomes; (ii) more mobile plasmids relative
95 to less mobile plasmids¹⁴⁻²⁰. We used phylogeny-based statistical methods that control for the
96 problem of non-independence, analysing 1632 genomes from 51 bacterial species, to examine
97 the location of genes that code for extracellular proteins. We then used theoretical models, to

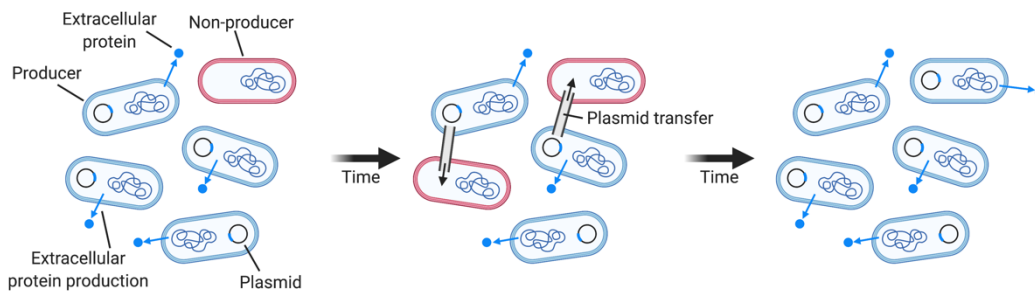
98 examine whether horizontal gene transfer facilitates the evolution as well as the initial spread
 99 of cooperation.

100

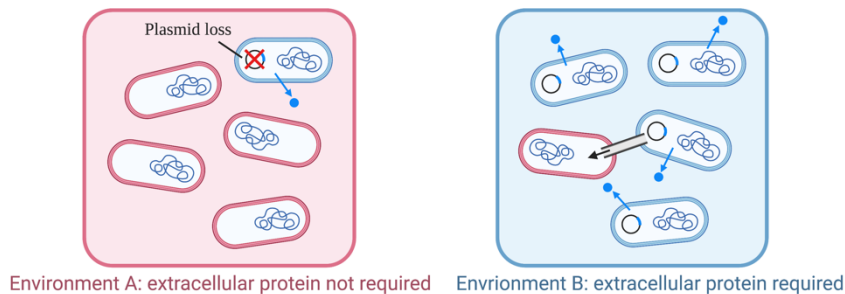
101 Finally, we also tested alternative hypotheses for why genes coding for extracellular proteins
 102 might be preferentially carried on plasmids. We used three measures of environmental
 103 variability to ask whether species which had more variable environments were those most
 104 likely to carry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

108

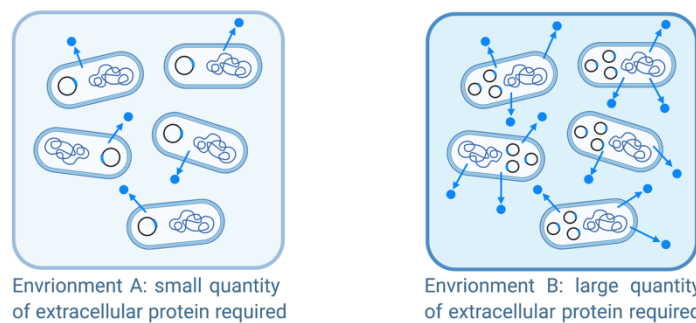
(a) Cooperation Hypothesis: Plasmid transfer stabilises cooperation by 'infecting' non-producing cheats



(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



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

111 (a) Cooperation Hypothesis. Blue cells produce extracellular proteins which act as cooperative
112 public goods, while red cells are ‘cheats’ which exploit this cooperation. Over time cheats grow
113 faster than cooperators since they forgo the cost of public good production. However, because
114 the gene for the extracellular protein is located on a plasmid, cooperators can transfer the gene
115 to the cheats, turning them into cooperators, increasing genetic relatedness at the cooperative
116 locus, and stabilising cooperation. (b) Gain and Loss Hypothesis. The production of the
117 extracellular protein is required in some environments, but not others. Transitions between
118 these environments can result from temporal or spatial change. Cells are selected to either lose
119 (Environment A) or gain (Environment B) the plasmid coding for the production of the
120 extracellular protein. (c) Beyond Horizontal Gene Transfer Hypothesis. The location of a gene
121 on a plasmid could provide a number of benefits, other than the possibility for horizontal gene
122 transfer³⁸. For example, when the quantity of extracellular protein required varies across
123 environments (A versus B), plasmid copy number could be varied to adjust production³⁸.

124

125 **Results**

126 **Genomic Analyses.**

127 We use the approach developed by Nogueira *et al.*^{15,19,20}, of using PSORTb³⁹ to predict the
128 subcellular location of every protein encoded by 1632 complete genomes from 51 diverse
129 bacterial species (Figure S1; Table S3). We are also building upon the work of researchers who
130 pointed out that extracellular (secreted) proteins are likely to provide a benefit to the local
131 population of cells, and hence act as cooperative public goods^{2,15,19,20,40}. The advantage of this
132 method is that it allows a large number of genes to be examined, across multiple species.

133

134 Overall, we found the average bacterial genome had 2696 protein-coding genes on the
135 chromosome(s), and 223 on the plasmid(s) (Table 1). Of these, an average of 57 genes (~2%)
136 coded for the production of an extracellular protein. These patterns are very similar to those
137 found previously^{15,19,20}. We followed methods from previous studies by assuming that genes
138 coding for extracellular proteins are more likely to represent public goods, since the diffusion
139 of these secreted proteins will often mean their effects are shared among neighbouring
140 cells^{3,15,19,20}.

141

142

| | Extracellular | Non-Extracellular | % Extracellular |
|---------------|---------------|-------------------|-----------------|
| Chromosome(s) | 52 | 2644 | 1.9% |
| Plasmid(s) | 5 | 218 | 2.4% |

143

144 **Table 1. Summary of location of genes coding for extracellular proteins across species.**

145 We calculated the mean number of genes coding for extracellular proteins and non-
146 extracellular proteins for all genomes in each species. We then calculated the mean of these
147 species means to give the values in the above table. The values above therefore provide an
148 indication of the location of genes coding for extracellular proteins in an average genome,
149 controlling for number of genomes per species. Genes with unknown protein localisations were
150 not included (Chromosome: 26.2%; Plasmid: 38.3%).

151

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

153 We found that extracellular proteins were not more likely to be carried on plasmids compared
154 to chromosomes (Figure 2). The difference in the proportion of genes that coded for
155 extracellular proteins between plasmid and chromosome was not significantly different from
156 zero across all species (MCMCglmm⁴¹; posterior mean = 0.004, 95% CI = -0.063 to 0.057,
157 pMCMC= 0.87; n = 1632 genomes; R² of species sample size = 0.47, R² of phylogeny = 0.17;
158 Table S2, row 1). This result was robust to alternative forms of analysis. We also found no
159 significant difference when we: (i) compared chromosomes to plasmids of only certain
160 mobilities (Fig S4; Table S2, rows 20-22); (ii) analysed our data by two alternative methods,
161 by looking at the ratio of proportions instead of the difference, or by considering only whether
162 the plasmid proportion was greater than the chromosome proportion, removing any effect of
163 the magnitude of this difference (Figure S5; Table S2, rows 2 and 3).

164

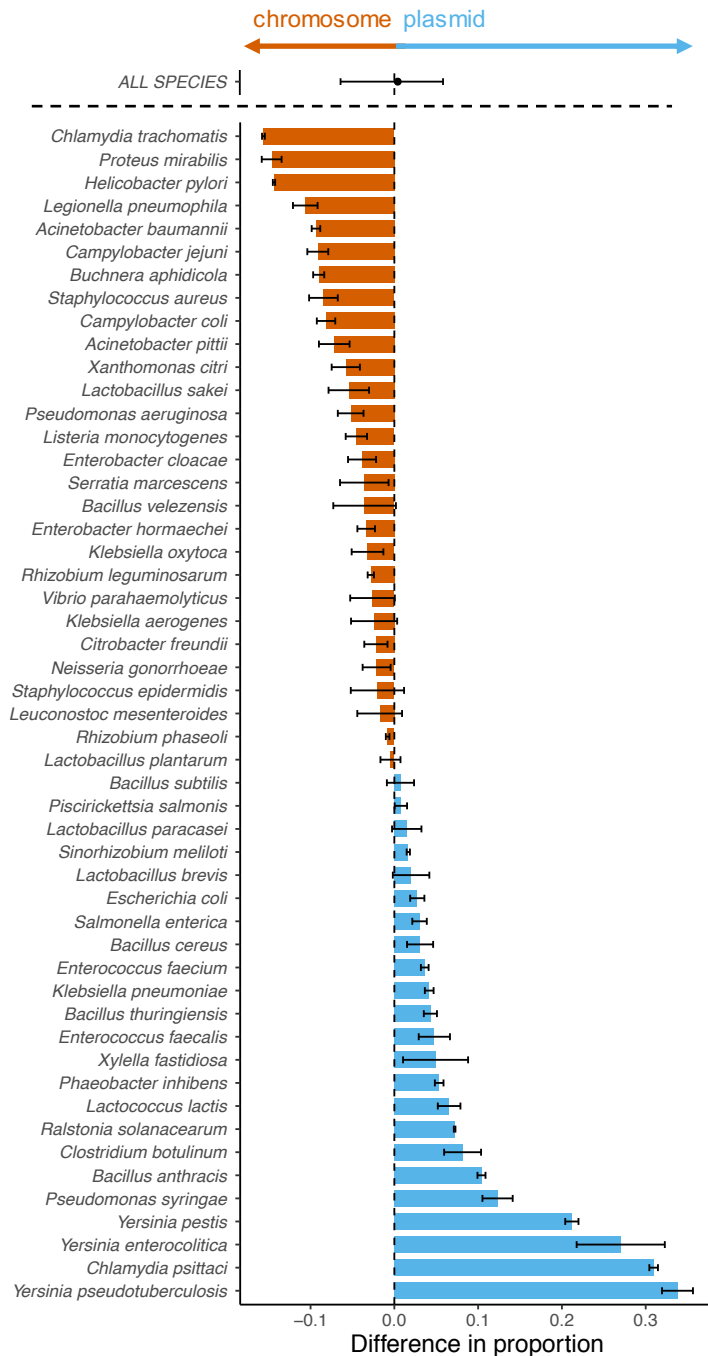
165 The lack of an overall significant result was clear when looking at the raw data for the different
166 species that we examined (Figure 2; Figure S5). There was considerable variation across
167 species in the location of genes coding for extracellular proteins. Overall, extracellular proteins
168 were more likely to be on plasmids in 51% of species (26/51), and more likely to be on the
169 chromosome(s) in 49% (25/51) of species (Figure S5). For example, in *Bacillus anthracis*
170 genes coding for extracellular proteins were three times more likely to be on plasmids, whereas
171 in *Acinetobacter baumannii* genes coding for extracellular proteins were three times more

172 likely to be on the chromosome(s) (Figure S5). Clearly, across species, genes coding for
173 extracellular proteins are not consistently more likely to be on plasmids.

174

175 As a control, we also analysed the genomic location of the genes coding for all other classes of
176 protein (Figure S1). Specifically, we analysed genes that coded for the production of
177 Cytoplasmic, Cytoplasmic Membrane, Periplasmic, Outer Membrane and Cell Wall proteins.

178 We found that none of these protein localisations were significantly overrepresented on
179 plasmids or chromosomes across the 51 species (Figure S6; Table S2, rows 5-10). Plasmids
180 are highly variable in the genes they carry.



181

182

183

184

185

186

187

188

189

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,

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

194

195 **Importance of controlling for non-independence of genomes.** Our results contrast with
196 previous studies, which found that plasmid genes code for proportionally more extracellular
197 proteins than chromosomes^{15,19,20}. The first of these studies found this pattern across 20
198 *Escherichia coli* genomes¹⁵. We also found that genes coding for extracellular proteins in *E.*
199 *coli* were more likely to be found on plasmids (Figure 2; Figure S5). However, Figure 2 shows
200 that this is not a consistent pattern across species: approximately half (25/51) of the species we
201 analysed showed a pattern in the opposite direction, with genes coding for extracellular proteins
202 more likely to be on their chromosome(s) than their plasmid(s).

203

204 Two subsequent, multi-species studies found that plasmid genes were significantly more likely
205 to code for extracellular proteins than chromosome genes^{19,20}. These studies used statistical
206 tests such as Wilcoxon signed-rank test to ask whether there was a consistent pattern, using
207 bacterial genomes as independent data points. When we analysed our data with the same
208 statistical methods used in these studies, we also obtained a significant result (Wilcoxon
209 signed-rank test; $V = 826530$, $p\text{-value} < 0.001$, $R^2 = 0.385$; $n = 1632$ plasmid-chromosome
210 pairs). When analysing other questions, Garcia-Garcera & Rocha²⁰ used MCMCglmm to
211 control for phylogeny.

212

213 Why does using bacterial genomes as independent data points lead to a significant result? By
214 using a Wilcoxon signed-rank test, at the level of the genome, we are implicitly assuming that
215 all the genomes analysed are: (i) independent from one another; (ii) a representative sample of
216 bacteria in nature. Neither of these are true for multi-species genomic datasets. First, due to
217 shared ancestry, species are not independent from one another, and so neither are genomes in
218 such analyses^{24,42}. Even a slight lack of independence can lead to heavily biased results in
219 statistical analyses and spurious conclusions²¹. Second, genomic databases tend to have a
220 disproportionate abundance of certain species and genera. This will bias the results towards
221 commonly sequenced species.

222

223 Consequently, when asking questions across species, it is inappropriate to treat all the genomes
224 in genomic datasets as independent data points. When we performed an analysis analogous to
225 the Wilcoxon signed-rank test, using the same untransformed data which produced a significant
226 result above, but controlled for the number of genomes per species and the non-independence
227 of species, we no longer found any significant difference between the proportion of plasmid
228 and chromosome genes coding for extracellular proteins (MCMCglmm; posterior mean =
229 0.017, 95% CI = -0.021 to 0.057, pMCMC = 0.332; n = 1632 plasmid-chromosome paired
230 differences in extracellular proportion; R^2 : species sample size = 0.46, phylogeny = 0.34; Table
231 S2, row 4). Furthermore, we found that the number of genomes per species and the non-
232 independence of species explained 46% and 34% of the variation in data respectively (paired
233 plasmid and chromosome differences across our 1632 genomes). Taken together, this
234 illustrates that it is not our data which disagrees with previous studies, but instead our use of
235 statistical analyses appropriate for multi-genome, multi-species datasets²³⁻²⁵.

236

237 These data also illustrate the importance of examining effect sizes, and not just whether results
238 are statistically significant. With large sample sizes it is possible to get results that are
239 significant but not biologically important. One rule of thumb is to assume that a result is only
240 biologically significant if the percentage of variance explained is >10% (i.e. $R^2 > 0.1$)⁴³. When
241 bacterial genomes are assumed to be independent data points in across species analyses, this
242 leads to inflated sample sizes. Consequently, even when results are statistically significant at
243 $P < 0.05$, they can still only explain 1-2% of the variation in the data, which is clearly not
244 biologically significant. The flip side of such considerations is that effects sizes and
245 examination of raw data at the species level (e.g. Figure 2) are also useful checks against non-
246 significant results due to a lack of statistical power (type II errors).

247

248 **Plasmids with higher mobility do not carry more genes for extracellular** 249 **proteins.**

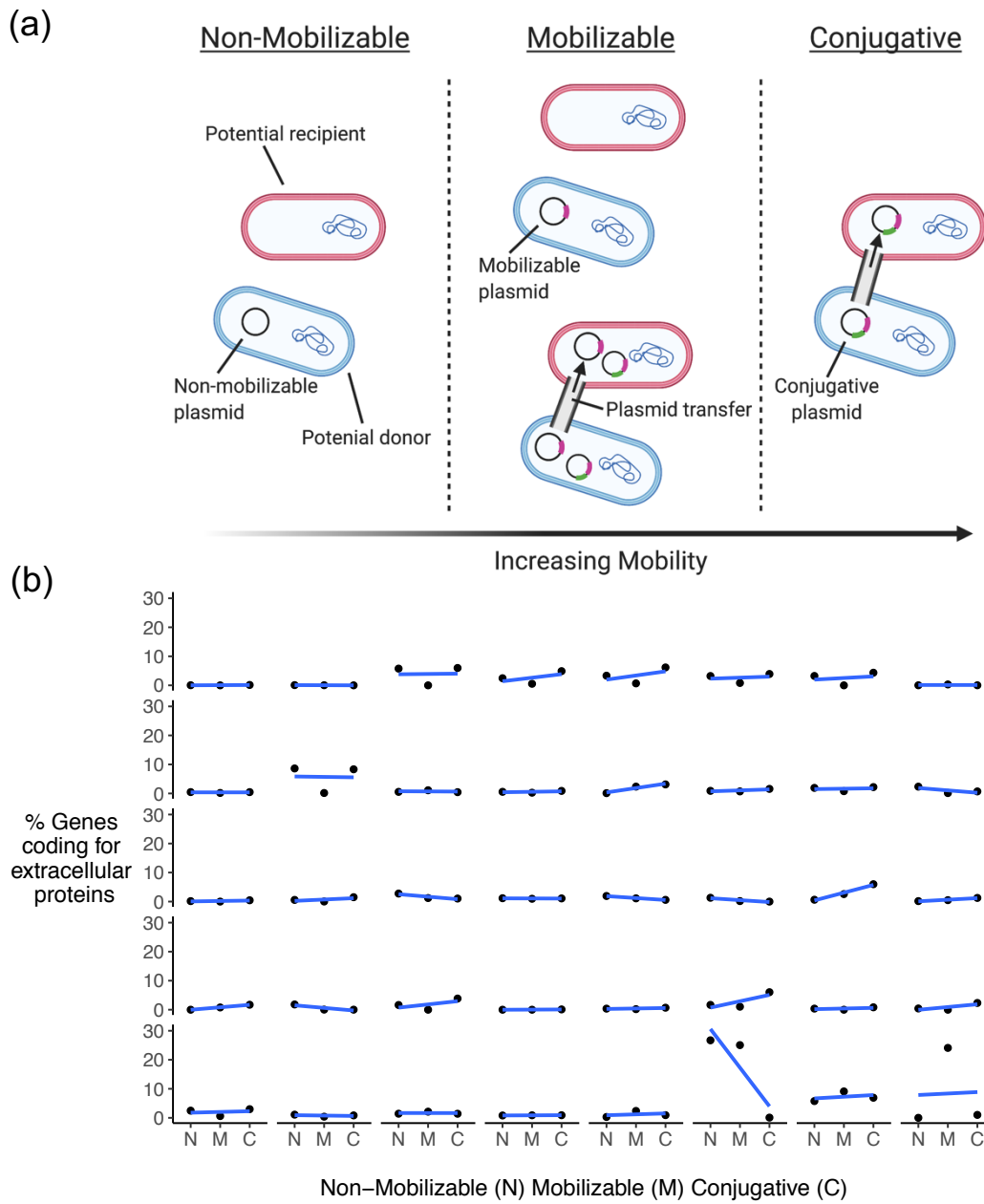
250 We then tested another prediction of the cooperation hypothesis: cooperation is more likely to
251 be favoured when coded for on more mobile plasmids¹⁴⁻¹⁸. We used data from the MOBsuite
252 database to assign plasmids to one of three levels of mobility (Fig 3a)^{44,45}. We classify:
253 conjugative plasmids, which carry all genes necessary to transfer, as the most mobile;
254 mobilizable plasmids, which are dependent upon conjugative plasmids' machinery to transfer,

255 to have intermediate mobility; non-mobilizable plasmids, which cannot be transferred via
256 conjugation, to be the least mobile (Fig 3a)^{44,46}.

257

258 Genes coding for extracellular proteins were not more likely to be on plasmids with higher
259 transfer rates (Figure 3b). Examining the slope of the regression between plasmid mobility and
260 the proportion of genes coding for extracellular proteins, we found no consistent pattern across
261 species (MCMCglmm; posterior mean = 0.006, 95% CI = -0.040 to 0.052, pMCMC = 0.73; n
262 = 40; Table S2, row 11). This lack of a significant relationship was robust to different forms of
263 analysis, including an examination of the means of each mobility type of each species (Figure
264 S7; Table S2, row 12). A caveat here is that our estimates of transfer rates across different types
265 of plasmid is relative, and it would be very useful to obtain quantitative estimates of transfer
266 rates.

267



269

270 **Figure 3. Plasmid mobility and extracellular proteins.** (a) We divided plasmids into three
 271 mobility types: non-mobilizable (lowest or no mobility); mobilizable (intermediate mobility);
 272 conjugative (highest mobility). Blue cells are potential plasmid donors, while red cells are
 273 potential recipients. Each panel shows when plasmid transfer is possible for one of the three
 274 plasmid mobility types. Non-mobilizable plasmids cannot be transferred. Mobilizable plasmids
 275 cannot be transferred alone, but they carry enough genes to ‘hijack’ the machinery of a
 276 conjugative plasmid that is in the same cell. Conjugative plasmids carry all genes necessary to
 277 transfer independently. (b) The 40 species which carried plasmids of all three mobilities are

278 shown, with a panel for each of these species. Dots in each panel indicate the mean % of genes
279 coding for extracellular proteins of all plasmids of each mobility level. The blue lines are the
280 linear regression of these three points. We arcsine square root transformed proportion data
281 before calculating the mean for each species, and then back-transformed these values for
282 display of the data. Overall, there is no consistent trend for genes that code for extracellular
283 proteins to be on more mobile plasmids.

284

285 **Theoretical Stability of Cooperation**

286 We examined whether cooperative genes should be overrepresented on plasmids, relative to
287 the chromosome. First, horizontal gene transfer on a plasmid could allow cooperation to be
288 favoured in conditions where it would otherwise not be favoured¹⁴⁻¹⁷. For example, because
289 plasmid transfer can turn non-cooperators in to cooperators, and increase relatedness at the loci
290 for cooperation¹⁷. Second, even if horizontal gene transfer did not increase the range of
291 biological scenarios (parameter space) where cooperation was favoured, there could be
292 selection for cooperation to be coded for on a plasmid, rather than a chromosome.

293

294 We followed the lifecycle assumptions of Mc Ginty *et al.*¹⁷, assuming that the population is
295 divided into patches, which are each founded by N independent cells (Supp. Info. 4). Cells
296 reproduce clonally until there are a large number of cells per patch, after which they disperse.
297 Cells can carry a plasmid which is transferred with probability β between paired cells, and
298 which is costly (C_C) to carry. Individuals with the gene for cooperation produce a public good,
299 at a cost C_G , which generates a benefit B that is shared between all members of the patch. The
300 gene for cooperation can be on the plasmid or chromosome.

301

302 Consistent with previous analyses, we found that horizontal gene transfer on a plasmid can
303 initially help cooperation invade (Figure 4). Horizontal gene transfer increased the frequency
304 of cooperation, by turning non-cooperators into cooperators, which also increases relatedness
305 at the cooperative locus^{14-18,47}.

306

307 In contrast, we found that transfer on a plasmid did not increase the range of parameter space
308 where cooperation was maintained at evolutionary equilibrium (Fig 4a) (Supp. Info. 4).
309 Specifically, cooperation was only favoured when $RB - C_G > 0$, where R is the genetic relatedness
310 at the chromosomal (individual) level ($R = 1/N$). Cooperation was therefore only favoured when

311 it provided a kin selected benefit at the level of the chromosome (individual), as predicted by
312 Hamilton's rule^{48,49}.

313

314 Our model therefore suggests that horizontal gene transfer can help cooperation initially
315 invade, but then does not help maintain cooperation in the long term. As a plasmid approaches
316 fixation, any benefit of horizontal gene transfer is lost. Consequently, competition between
317 plasmids with and without a cooperative gene (cooperators and cheats) becomes analogous to
318 the scenario in which the gene for cooperation is on the chromosome. An analogous result was
319 also found in a meta-population model by Mc Ginty *et al.*¹⁶. Our prediction has been supported
320 experimentally by Bakkeren *et al.*³⁰, who found that location on a conjugative plasmid could
321 help a cooperative trait invade in *Salmonella* Typhimurium (*S.Tm*), but that this was only stable
322 with strong population bottlenecks (high relatedness).

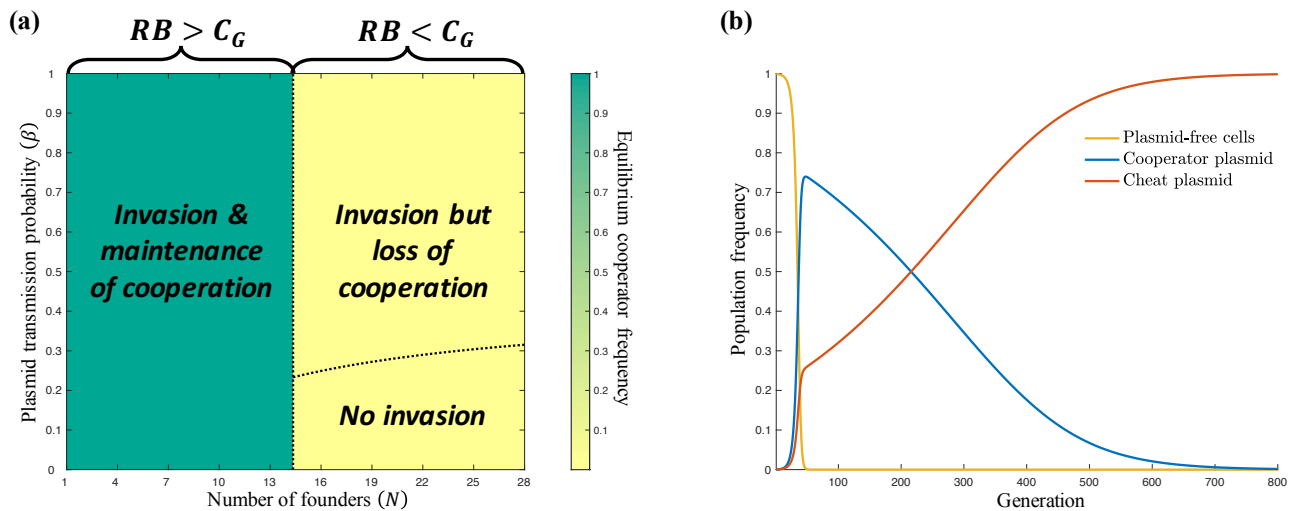
323

324 In addition, we found that, when cooperation is favoured, cooperative traits are not more likely
325 to be favoured on, or transferred to, plasmids. The reason is that, when cooperation is favoured,
326 non-cooperators (cheats) are purged from the population, which means there is no extra fitness
327 benefit of coding for the cooperative trait on a plasmid rather than the chromosome.
328 Consequently, our results suggest that horizontal gene transfer doesn't necessarily favour
329 cooperation. Our results differ from previous theory because we have examined both: (i) a
330 greater range of genetic architectures, especially plasmids that do not encode cooperation; and
331 (ii) the evolutionary stability (maintenance) of cooperation, not just its initial invasion, while
332 explicitly modelling plasmid transmission (Supp. Info. 4)^{16,47}.

333

334 More generally, with regard to whether we should expect plasmids and chromosomes to be in
335 conflict, our results emphasise that at evolutionary equilibrium, the fitness interests of plasmids
336 and chromosomes can be expected to align. Although, there could be interesting transient
337 dynamics, where conflict leads to cooperation being favoured temporarily (Figure 4b). Another
338 important factor is the rate of horizontal gene transfer. While plasmids clearly transmit fast
339 enough to influence evolution, the transfer rates per cell per generation do not appear high
340 enough to significantly influence relatedness at the locus for cooperation⁵⁰.

341



342

343 **Figure 4. Plasmids facilitate the invasion but not the maintenance of cooperation.** In parts
344 (a) and (b), we plot the results of our theoretical model. (a) Cooperation is only maintained at
345 equilibrium (green shaded area) when it is favoured at the chromosomal level $RB > C_G$, which
346 is unaffected by plasmid transfer (β). (b) Plasmids can facilitate the invasion and initial spread
347 of cooperation (blue line shoots above red line), but cooperative plasmids are eventually
348 outcompeted by cheat plasmids (red line goes to 1). To generate the plots in (a) and (b), we
349 assumed the following parameter values: (a & b) $B = 1.435, C_G = 0.1, C_C = 0.2$; (b) $\beta =$
350 $0.5, N = 16$.

351

352 **Alternate hypotheses.**

353 Finally, we examined whether alternate hypotheses may better explain the considerable
354 variation in the location of genes coding for extracellular proteins across species. Species which
355 live in more variable environments may be more likely to carry extracellular genes on plasmids.
356 This could be expected for different reasons, including plasmid transfer allowing genes for
357 different environments to be gained and lost (Figure 1b), or plasmids conferring some other
358 advantage not associated with horizontal gene transfer, such as allowing copy number to be
359 conditionally adjusted (Figure 1c)^{31,32,38,51}. There are a number of different ways to classify
360 environmental variability, and so we used three different methods.

361

362 **Broad host-range pathogens are most likely to carry genes for extracellular proteins on**
363 **plasmids.** We first used the diversity of pathogen hosts as a proxy for environmental
364 variability. Although this does not capture all environmental variability experienced by species

365 in our data set, pathogenicity is a key aspect of bacterial lifestyle that has been suggested to be
366 important for plasmid gene content, such as antibiotic resistance and virulence factors^{6,40,52,53}.
367 We divided species into three categories: pathogens with broad host-range, pathogens with
368 narrow host-range, and non-pathogens. Broad host-range pathogens are expected to encounter
369 more variable environments than narrow host-range pathogens.

370

371 We found that pathogens with a broad host-range were more likely to carry genes coding for
372 extracellular proteins on their plasmids, compared with both narrow host-range pathogens and
373 non-pathogens (Fig 5). Specifically, we compared the difference in the proportion of genes
374 coding for extracellular proteins between plasmid(s) and chromosome(s) across these three
375 categories of species (MCMCglmm; Narrow compared to Broad host-range pathogens:
376 posterior mean = -0.222, 95% CI = -0.322 to -0.123, pMCMC = <0.001; Non-pathogens
377 compared to Broad host-range pathogens: posterior mean = -0.161, 95% CI = -0.252 to -0.067,
378 pMCMC = <0.001; n = 701 genomes; R² of pathogenicity/host-range = 0.35, R² of species
379 sample size = 0.28, R² of phylogeny = 0.11; Table S2, row 23). There was no significant
380 difference between narrow host-range pathogens and non-pathogens in the proportion of genes
381 coding for extracellular proteins on their plasmids compared to chromosome(s) (MCMCglmm;
382 Non-pathogens compared to Narrow host-range pathogens: posterior mean = 0.031, 95% CI =
383 -0.065 to 0.127, pMCMC = 0.482; n = 389; Table S2, row 25). These patterns hold irrespective
384 of whether we included species that we could not reliably classify into either category, such as
385 opportunistic pathogens, in our analyses (Figure S10).

386

387 **Plasmids of broad host-range pathogens carry many pathogenicity genes.** We suspected
388 that the additional extracellular proteins coded for by plasmids of broad host-range species,
389 compared to narrow host-range species, may be particularly involved in facilitating
390 pathogenicity^{40,52,53}. To investigate this, we used the program MP3⁵⁴ to assign a each
391 extracellular protein as either 'pathogenic' or 'non-pathogenic'.

392

393 We found that plasmids of broad host-range pathogens were particularly enriched with
394 extracellular proteins involved in facilitating pathogenicity, compared to plasmids of narrow
395 host-range species (Figure 6). Specifically, we found that pathogens with a broad host-range
396 were significantly more likely to code for pathogenic extracellular proteins on their plasmids
397 compared to narrow host-range species (Figure 6a) (MCMCglmm; Narrow compared to Broad
398 host-range pathogens: posterior mean = -0.209, 95% CI = -0.350 to -0.086, pMCMC = 0.012;

399 n=474 genomes; Table S2, row 26). In contrast, the relative location of non-pathogenic
400 extracellular proteins did not vary between broad and narrow host-range pathogens (Figure 6b)
401 (MCMCglmm; Narrow compared to Broad host-range pathogens: posterior mean = -0.036,
402 95% CI = -0.115 to 0.040, pMCMC = 0.296; n=474 genomes; Table S2, row 27).
403 Consequently, the excess of genes coding for extracellular proteins on the plasmids of broad
404 host-range species (Figure 5) appears to arise due to an excess of pathogenicity genes coding
405 for extracellular proteins (Figure 6).

406

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

415

416 Pathogenic extracellular proteins could be preferentially coded for on plasmids to facilitate
417 their gain and loss (Figure 1b: Gain and loss hypothesis), or because of some other benefit
418 provided by being carried on a plasmid (Figure 1c: Beyond horizontal gene transfer
419 hypothesis). We tested these possibilities by examining whether pathogenic extracellular
420 proteins were more likely to be on plasmids that transfer at higher rates. This would be
421 predicted by the gain and loss hypothesis, but not the beyond horizontal gene transfer
422 hypothesis. We found that plasmids with higher mobility did not code for more pathogenic
423 extracellular proteins. Specifically, across broad host-range pathogen species, the slope of the
424 regression between plasmid mobility and the proportion of genes coding for pathogenic
425 extracellular proteins was not consistently positive (Figure S11) (MCMCglmm; posterior mean
426 = -0.020, 95% CI = -0.224 to 0.185, pMCMC = 0.774; n=7; Table S2, row 31). This lack of a
427 significant relationship was robust to additional forms of analysis, such as considering all
428 pathogenic species, including narrow host-range pathogens and those not carrying plasmids of
429 all three mobility types (Figure S12; Table S2, rows 32 and 33).

430

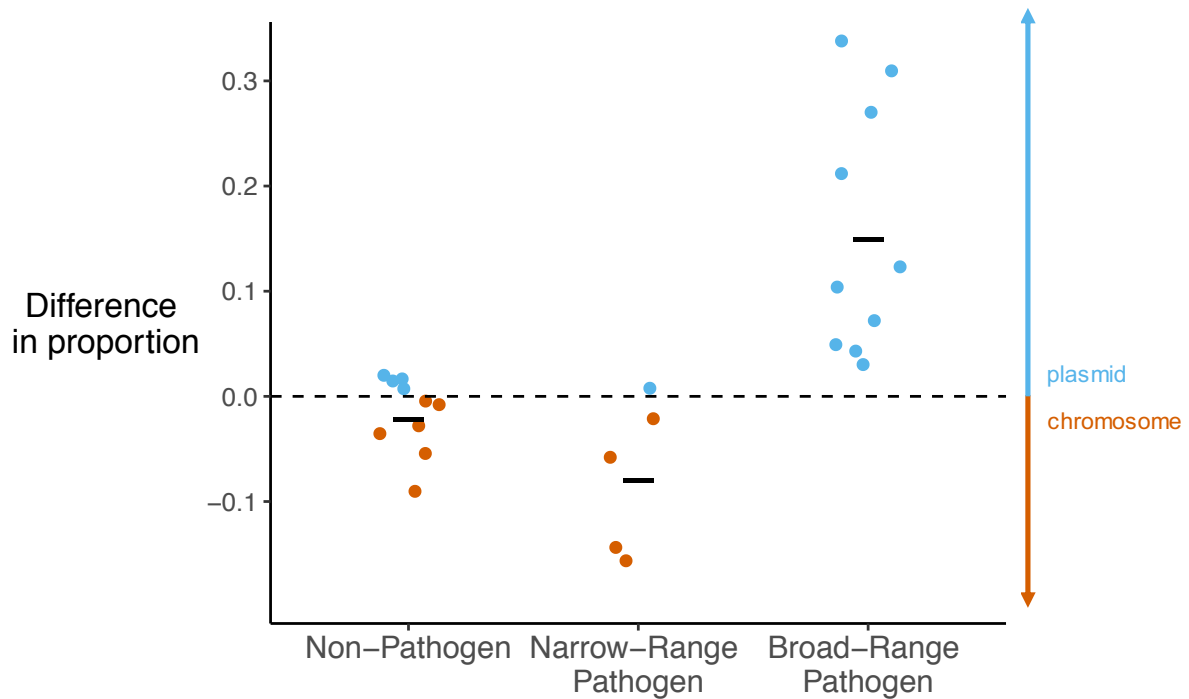
431 Taken together, our results are most consistent with the hypothesis that genes coding for
432 extracellular proteins are overrepresented on plasmids when plasmid carriage provides a

433 benefit other than mobility (Figure 1c). A number of other factors may influence which genes
434 are carried on plasmids, beyond horizontal gene transfer. First, there is evidence that increasing
435 the copy number of plasmids can lead to increasing rates of evolution in the genes they carry⁵⁵,
436 and it also may act as a mechanism to increase the expression of genes carried on plasmids^{56,57}.
437 For example, increased expression of genes coding for extracellular public goods such as
438 virulence factors could help invasion of a host and utilisation of host resources. This could be
439 particularly beneficial for broad host-range pathogens that frequently invade a variety of
440 different hosts. Copy number of plasmids has also recently been shown to lead to genetic
441 dominance effects⁵¹, with likely implications for the phenotypes of genes selected for plasmid
442 carriage⁵¹. Second, plasmids compete with their bacterial hosts for resources such as replication
443 machinery and nucleotides^{58,59}. To resolve this competition, plasmids should be under selection
444 to reduce their cost to the host, with a likely impact on their gene content. For example,
445 extracellular proteins are, on average, cheaper to produce than intracellular proteins^{15,20}.
446 Plasmid-host competition could consequently select for plasmids to carry more genes coding
447 for cheaper proteins, and so more extracellular proteins. Taken together, there are a number of
448 factors which will allow plasmids to facilitate pathogenicity and adaptation to new and variable
449 environments. Our conclusion here should be seen as tentative, as some form of the gain and
450 loss hypothesis (Figure 1b) could still be argued to be consistent with the data, if it is just the
451 potential for horizontal gene transfer that matters, and not the rate.

452

453

454



455

456 **Figure 5. Environmental variability and the location of genes coding for extracellular**

457 **proteins.** We have divided species into either pathogens or non-pathogens, with pathogens

458 further categorised into those with a narrow or broad host-range. The y-axis shows the

459 difference in the proportion of genes on plasmids and chromosomes coding for extracellular

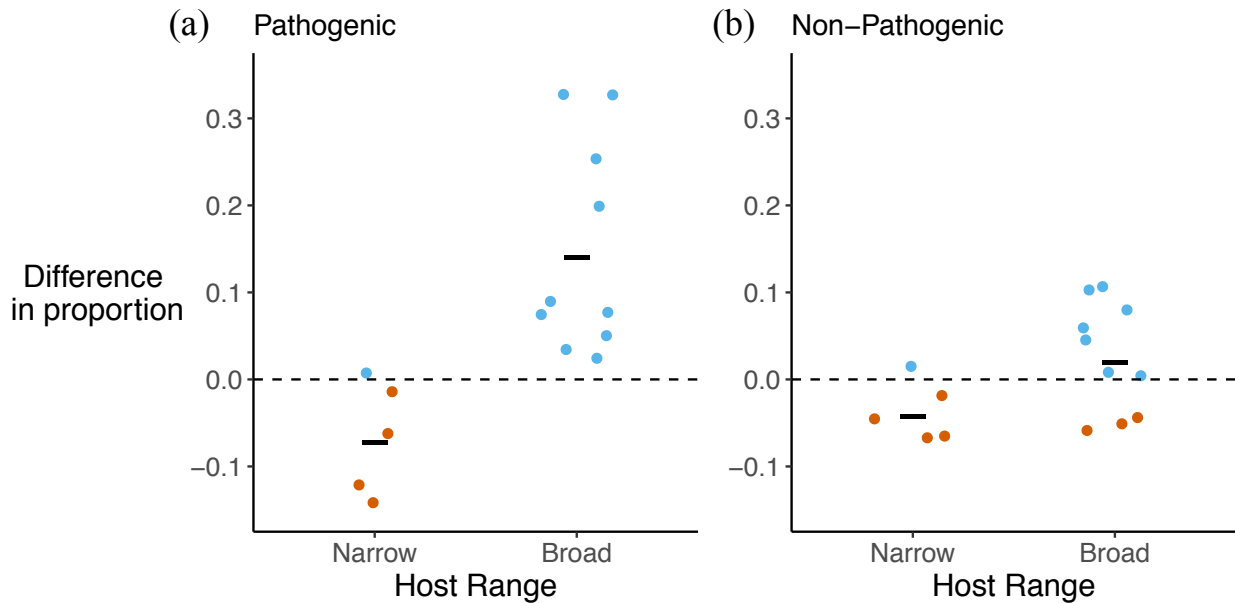
460 proteins. Each dot is the mean for all genomes in a species. Species in blue are those with

461 extracellular proteins overrepresented on plasmids, while species in red are those with

462 extracellular proteins overrepresented on chromosomes. The black bars indicate the mean for

463 all species in each category. Overall, pathogens with a broad host-range are more likely to have

464 genes coding for extracellular proteins on their plasmids.



465 **Figure 6. The location of genes coding for pathogenic and non-pathogenic extracellular**
 466 **proteins, in species with broad and narrow host-ranges.** We categorised pathogenic species
 467 into those with either a broad or narrow host-range. The y-axes in (a) and (b) show the
 468 difference in the proportion of genes coding for extracellular proteins on plasmids and
 469 chromosomes which are predicted by MP3 as either (a) pathogenic or (b) non-pathogenic.
 470 Higher values indicate that extracellular proteins are more likely to be coded for by plasmids.
 471 Each dot is the mean for all genomes in a species. Species in blue are those with the relevant
 472 subset of extracellular proteins overrepresented on plasmids, while species in red are those with
 473 the subset of extracellular proteins overrepresented on chromosomes. Overall, there is a
 474 significant difference between broad and narrow host-range species in the location of genes
 475 coding for pathogenic extracellular proteins, but no difference for non-pathogenic extracellular
 476 proteins.

477

478 **Number of environments and core vs accessory genes.** We also looked at two additional
 479 measures of environmental variability: (i) the number of five broad environments a species was
 480 sequenced in^{20,60,61}; (ii) the proportion of a species' genomes that is composed of 'core' genes,
 481 which are those found in all genomes of the species – species which experience more variable
 482 environments appear to have relatively smaller core genomes³². We found no significant
 483 correlation between either of these measures and the likelihood that genes coding for
 484 extracellular proteins were carried on plasmids (Figure S13) (Supp. Info. 1; Table S2, rows 35
 485 and 37). Garcia-Garcera & Rocha²⁰ previously analysed a different but related question,

486 examining the type of environment, and also used a MCMCglmm to control for the
487 phylogenetic structure of the data (Supp. Info. 1).

488

489 **Complementary Analyses**

490 There a number of directions in which our analyses could be expanded. We focused on
491 plasmids because they have been the focus of previous theoretical and empirical work^{14,16-18}.
492 Other mobile genetic elements include bacteriophages and integrative conjugative
493 elements^{62,63}. Comparing core and accessory genes could be a potential way to lump all causes
494 of horizontal gene transfer^{15,19}. We considered the relative transfer rates among mobility types;
495 quantitative estimates of plasmid transfer rates would be very useful for further examination of
496 plasmid mobility^{46,50,64-66}. We followed previous genomic studies by using extracellular
497 proteins as indicators of cooperative traits^{2,15,19,20}. The advantages of this approach are that: (i)
498 we could compare our results with those from previous studies; (ii) secretion systems are highly
499 conserved, allowing us to examine a large number of species, where detailed genetic
500 annotations are lacking; (iii) cooperation mediated by extracellular proteins is usually
501 controlled by only one gene, making them potentially more suitable for plasmid carriage
502 compared to cassettes of multiple genes^{67,68}. However, while extracellular proteins are likely
503 to be cooperative traits, not all cooperative genes code for extracellular proteins (e.g. secondary
504 metabolites such as siderophores), and not all extracellular proteins are involved in cooperation
505 (e.g. those involved in motility such as flagellin). It would be very useful to examine more
506 detailed annotations of social genes, and expand to other mobile genetic elements.

507

508 **Discussion**

509 We found no support for the hypothesis that horizontal gene transfer favours cooperation. Our
510 genomic analyses showed that extracellular proteins are not: (i) overrepresented on plasmids
511 compared to chromosomes; (ii) more likely to be carried by plasmids that transfer at higher
512 rates. These patterns could be explained by theoretical modelling, which showed that while
513 horizontal gene transfer may help cooperation to initially invade a population, it does not then
514 help the maintenance of cooperation in the long term. Once plasmids become common, cheat
515 plasmids that do not code for cooperation are able to outcompete cooperative plasmids,
516 analogous to selection at the level of the chromosome¹⁶. Our prediction has also been supported
517 experimentally by Bakkeren *et al.*³⁰, in *Salmonella* Typhimurium (*S.Tm*), who observed
518 cooperation invading on a plasmid, but then being outcompeted by newly emerging non-

519 cooperative cheats. In contrast, we found that genes coding for extracellular proteins involved
520 in pathogenicity and virulence are preferentially located on plasmids in pathogens with a broad
521 host-range. These pathogenic virulence genes were not preferentially located on plasmids that
522 transfer at a higher rate, suggesting that the benefit of being located on a plasmid is something
523 other than horizontal gene transfer, such as the ability to vary copy number.

524

525 **Methods**

526 **Genome Collection**

527 We retrieved 1632 complete genomes comprising 51 bacterial species from GenBank RefSeq
528 (<https://www.ncbi.nlm.nih.gov>) between February-November 2019. We used species on panX
529 (<http://pangenome.tuebingen.mpg.de>)⁶⁹ as a list of potential species for our dataset, since these
530 comprise the most sequenced bacterial species. To allow comparison of chromosome and
531 plasmid genes within the same genome, we only retrieved genomes that contained at least one
532 plasmid sequence. We included species with 10 or more RefSeq genomes with one or more
533 plasmids available in our analysis. We retrieved up to 100 genomes for each species; this was
534 either all complete genomes available for the species, or a random sample where more than
535 100 were available. Where two or more genomes had the same strain name, we randomly
536 retrieved one genome to reduce the risk of pseudoreplication.

537

538 **Prediction of Subcellular Location of Proteins**

539 We used PSORTb v.3³⁹ to predict the subcellular location of every protein encoded by each
540 genome in our dataset. We used a Docker image of PSORTb developed by the Brinkman Lab,
541 available at: https://github.com/brinkmanlab/psortb_commandline_docker. We chose
542 PSORTb because it is widely regarded as one of the best performing programs of its kind⁷⁰. It
543 has also been used in previous analyses to identify ‘cooperative’ genes and/or extracellular
544 proteins in bacteria^{15,20}. The program has a number of modules which are trained to recognise
545 particular features of proteins. Results from these modules are combined to give a Final
546 Prediction for each protein. We consulted the literature to confirm the Gram stain of each of
547 our species. For Gram-positive species, PSORTb assigns proteins to one of four locations
548 within the cell: cytoplasmic, cytoplasmic membrane, extracellular or cell wall (Figure S1). The
549 locations for Gram-negative species are the same, except that cell wall is replaced with outer
550 membrane and periplasmic, meaning there are five possible locations for proteins of Gram-
551 negative species (Figure S1). We used these predicted locations throughout all subsequent

552 analyses in this work. PSORTb could not reliably assign a subcellular location to 27% of
553 proteins we analysed, giving a final prediction of ‘unknown’ (Table S1). Unless explicitly
554 stated, we did not include these unknown proteins in our analyses.

555

556 **Predicting Plasmid Mobility**

557 We also predicted the mobility of every plasmid in our dataset using the MOB-typer tool of
558 the program MOBsuite⁴⁴. This searches for features of plasmid sequences including the origin
559 of transfer (oriT), relaxase and mating-pair formation to give each plasmid one of three
560 mobility predictions: (i) conjugative, where plasmids encode all machinery required to transfer
561 via conjugation; (ii) mobilizable, where plasmids do not encode all machinery, but encode oriT
562 and/or relaxase, allowing them to ‘hijack’ another plasmid’s conjugation machinery and
563 mobilize; (iii) non-mobilizable, where plasmids do not encode the genes necessary to be
564 mobilized by themselves or other plasmids, and so cannot transfer via conjugation. 628 of the
565 4150 plasmids in our dataset were flagged as ‘unverified’ against the MOBsuite dataset,
566 meaning their mobility prediction was unreliable and they were not included. This left 3522
567 plasmids for subsequent analysis.

568

569 **Effect of Mobility on Plasmid Extracellular Protein Content**

570 We next examined how plasmid mobility correlates with each plasmid’s extracellular protein
571 proportion. As part of its mobility prediction, MOBsuite⁴⁴ identifies sequences within each
572 plasmid involved with conjugation. To control for the possibility that conjugative plasmids, by
573 definition of being conjugative, must carry genes controlling this process, we subtracted the
574 total number of these sequences from the total number of proteins when calculating the
575 extracellular proportion of each plasmid. This is a highly conservative control, since it assumes
576 none of the proteins predicted as extracellular are involved in conjugation. We did all analyses
577 on these data with and without removing these mating-pair accessions to ensure any results
578 were not affected by factors unrelated to plasmids’ extracellular protein content.

579

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

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

586

587 **Measures of Bacterial Lifestyle and Environmental Variability**

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

599

600 To estimate the host-range of pathogens, we used information from the literature to determine
601 the maximum taxonomic level of hosts each species is able to invade. We defined narrow host-
602 range species as those which can invade either only one host species, or host species within the
603 same genus or family. In contrast, we defined broad-host range pathogens as those capable of
604 invading host species within the same order, class or phylum. For example, *Xanthomonas citri*
605 acts as a plant pathogen within the genus *Citrus*⁷¹, while *Pseudomonas syringae* acts as plant
606 pathogen across multiple orders of flowering plants⁷². For more details and references to the
607 literature used for this classification, please see Table S3.

608

609 We completed additional analyses for other two measures and proxies of environmental
610 variability, the details and results of which can be found in Supp. Info. 1. In brief, we used
611 previously published data which classified the habitat diversity of species using 16S RNA
612 environmental datasets across five broad habitats: water, wastewater, sediment, soil and
613 host^{60,61}. We also supplemented this with information from the literature for species not
614 included in the published data. We used this to ask whether species which lived in multiple
615 habitats had genes encoding extracellular proteins more overrepresented on their plasmids.

616

617 We also looked at bacterial pangenomes as a proxy for environmental variability, since it has
618 been noted that species with a high % of accessory genes, defined as genes found in only a
619 subset of genomes within a species, are generally those with more variable environments. All
620 pangenome data was collected from panX⁶⁹ (<http://pangenome.tuebingen.mpg.de>), since this
621 calculates the pangenome using the same method across all of our species.

622

623 **Pathogenicity categorisation of extracellular proteins**

624 We used MP3⁵⁴ to examine the pathogenicity of extracellular protein-coding genes in broad
625 host-range and narrow host-range pathogens. MP3 uses two modules to produce a ‘Hybrid’
626 prediction for each protein: either ‘Pathogenic’ or ‘Non-Pathogenic’. We used MP3 with
627 default parameters to gain this prediction for every extracellular protein in all genomes of broad
628 and narrow host-range species. MP3 was unable to give a prediction for approximately 9% of
629 extracellular proteins, and so these were not included in this analysis.

630

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

635

636 **Statistical analyses**

637 **MCMCglmm.** Many commonly used statistical methods in biology require data points to be
638 independent from one another. However, due to shared ancestry, species cannot be considered
639 as independent data points²⁴. Recently developed statistical methods now allow for
640 phylogenetic relationships to be controlled for within mixed effects models. For all statistical
641 analyses we used the MCMCglmm (Markov Chain Monte Carlo generalised linear mixed
642 effects model) package in R with phylogeny and sample size as random effects^{41,73}. We
643 extracted from each model the posterior mean, 95% Credible Intervals (functionally similar to
644 95% Confidence Intervals), and the pMCMC value (generally interpreted in a similar way to a
645 ‘p-value’). We also calculated R² values for models of particular interest using methods
646 described in^{74,75}. A detailed description of MCMCglmm can be found elsewhere^{41,73}.

647

648 The response variable in all of our analyses is either a proportion or a measure calculated from
649 proportions. Proportion data is bound between 0 and 1 and has a non-normal distribution. To

650 control for this, all proportion data in our analyses has been arcsine square root transformed to
651 improve normality.

652

653 **Phylogeny.** To control for species relationships, we generated a phylogeny including all 51
654 species in our dataset (Fig S2). We used a recently published maximum likelihood tree using
655 16S ribosomal protein data as the basis for our phylogeny⁷⁶. This tree of life typically had only
656 one representative species per genus. We used the R package ‘ape’ to extract all branches
657 matching species in our dataset⁷⁷. In cases where the genus representative was different to the
658 species in our dataset, we swapped the tip name with our species, since all members of the
659 same genus are equally related to members of a sister genus. In cases where we had multiple
660 species within a single genus in our dataset, we used the R package ‘phylotools’ to add these
661 species as additional branches into their genus⁷⁸. We used published phylogenies from the
662 literature to add any within-genus clustering of species’ branches. We used this phylogeny in
663 nexus format for all our MCMCglmm analyses (Fig S2, Table S2). Methods are also available
664 to control for uncertainty in phylogenetic reconstruction^{79,80}, although we have not done this
665 here.

666

667 **Acknowledgements**

668 We thank: Craig MacLean, Kevin Foster, Laurence Belcher, Chunhui Hao, and especially
669 Eduardo Rocha for their helpful comments; James Robertson for providing plasmid mobility
670 data from the MOBSuite database; the BBSRC (A.E.D.), ERC (J.L.T., T.W.S., A.S.G., M.G.
671 and S.A.W.), and NSERC-CRSNG of Canada (G.W.) for funding. Conceptual figures were
672 created with Biorender.com.

673

674 **Author Contributions**

675 A.E.D., J.L.T., A.S.G., S.A.W and M.G. conceived the genomic analyses and interpreted
676 results. A.E.D. and J.L.T. collected and analysed genomic data, and A.E.D. produced the
677 corresponding figures. T.W.S, G.W. and S.A.W. conceived the theoretical modelling and
678 interpreted results. T.W.S. completed the formal theoretical modelling. A.E.D., J.L.T, T.W.S.,
679 S.A.W., and M.G. wrote and/or edited the manuscript. A.E.D. wrote and put together S1, S2
680 and S3, and T.W.S. wrote and put together S4. All authors commented on and approved the
681 manuscript for submission.

682

683

684 **Competing Interests**

685 The authors declare no competing interests.

686

687 **Data Availability Statement**

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

691

692 **References**

- 693 1. Foster, K. R. Social behaviour in microorganisms. in *Social Behaviour* (eds. Szekely, T.,
694 Moore, A. J. & Komdeur, J.) 331–356 (Cambridge University Press, 2010).
695 doi:10.1017/CBO9780511781360.027.
- 696 2. McNally, L., Viana, M. & Brown, S. P. Cooperative secretions facilitate host range
697 expansion in bacteria. *Nat. Commun.* **5**, (2014).
- 698 3. West, S. A., Griffin, A. S., Gardner, A. & Diggle, S. P. Social evolution theory for
699 microorganisms. *Nat. Rev. Microbiol.* **4**, 597–607 (2006).
- 700 4. Simonet, C. & McNally, L. Kin selection explains the evolution of cooperation in the gut
701 microbiota. *Proc. Natl. Acad. Sci.* **118**, (2021).
- 702 5. Griffin, A. S., West, S. A. & Buckling, A. Cooperation and competition in pathogenic
703 bacteria. *Nature* **430**, 1024–1027 (2004).
- 704 6. Hale, T. L. Genetic basis of virulence in *Shigella* species. *Microbiol. Rev.* **55**, 206–224
705 (1991).
- 706 7. Dinges, M. M., Orwin, P. M. & Schlievert, P. M. Exotoxins of *Staphylococcus aureus*.
707 *Clin. Microbiol. Rev.* **13**, 16–34, table of contents (2000).
- 708 8. Diggle, S. P., Griffin, A. S., Campbell, G. S. & West, S. A. Cooperation and conflict in
709 quorum-sensing bacterial populations. *Nature* **450**, 411–414 (2007).
- 710 9. Jones, S. *et al.* The lux autoinducer regulates the production of exoenzyme virulence
711 determinants in *Erwinia carotovora* and *Pseudomonas aeruginosa*. *EMBO J.* **12**, 2477–
712 2482 (1993).
- 713 10. Sandoz, K. M., Mitzimberg, S. M. & Schuster, M. Social cheating in *Pseudomonas*
714 *aeruginosa* quorum sensing. *Proc. Natl. Acad. Sci.* **104**, 15876–15881 (2007).

- 715 11. Ghoul, M., Griffin, A. S. & West, S. A. Toward an evolutionary definition of cheating.
716 *Evolution* **68**, 318–331 (2014).
- 717 12. Butaitė, E., Baumgartner, M., Wyder, S. & Kümmerli, R. Siderophore cheating and
718 cheating resistance shape competition for iron in soil and freshwater *Pseudomonas*
719 communities. *Nat. Commun.* **8**, 414 (2017).
- 720 13. Thomas, C. & Nielsen, K. Thomas CM, Nielsen KM. Mechanisms of, and barriers to,
721 horizontal gene transfer between bacteria. *Nat Rev Micro* 3: 711-721. *Nat. Rev.*
722 *Microbiol.* **3**, 711–21 (2005).
- 723 14. Smith, J. The social evolution of bacterial pathogenesis. *Proc. R. Soc. Lond. B Biol. Sci.*
724 **268**, 61–69 (2001).
- 725 15. Nogueira, T. *et al.* Horizontal Gene Transfer of the Secretome Drives the Evolution of
726 Bacterial Cooperation and Virulence. *Curr. Biol.* **19**, 1683–1691 (2009).
- 727 16. Mc Ginty, S. E., Rankin, D. J. & Brown, S. P. Horizontal gene transfer and the evolution
728 of bacterial cooperation: mobile elements and bacterial cooperation. *Evolution* **65**, 21–32
729 (2011).
- 730 17. Mc Ginty, S. É., Lehmann, L., Brown, S. P. & Rankin, D. J. The interplay between
731 relatedness and horizontal gene transfer drives the evolution of plasmid-carried public
732 goods. *Proc. R. Soc. B Biol. Sci.* **280**, 20130400 (2013).
- 733 18. Dimitriu, T. *et al.* Genetic information transfer promotes cooperation in bacteria. *Proc.*
734 *Natl. Acad. Sci.* **111**, 11103–11108 (2014).
- 735 19. Nogueira, T., Touchon, M. & Rocha, E. P. C. Rapid Evolution of the Sequences and
736 Gene Repertoires of Secreted Proteins in Bacteria. *PLoS ONE* **7**, e49403 (2012).
- 737 20. Garcia-Garcera, M. & Rocha, E. P. C. Community diversity and habitat structure shape
738 the repertoire of extracellular proteins in bacteria. *Nat. Commun.* **11**, 758 (2020).
- 739 21. Kruskal, W. Miracles and Statistics: The Casual Assumption of Independence. *J. Am.*
740 *Stat. Assoc.* **83**, 929–940 (1988).
- 741 22. Ives, A. R. & Zhu, J. Statistics for correlated data: phylogenies, space, and time. *Ecol.*
742 *Appl. Publ. Ecol. Soc. Am.* **16**, 20–32 (2006).
- 743 23. Felsenstein, J. Phylogenies and the Comparative Method. *Am. Nat.* **125**, 1–15 (1985).
- 744 24. Harvey, P. H. & Pagel, M. D. *The Comparative Method in Evolutionary Biology.* (Oxford
745 University Press, 1991).
- 746 25. Grafen, A. The phylogenetic regression. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **326**,
747 119–157 (1989).

- 748 26. Hurlbert, S. H. Pseudoreplication and the Design of Ecological Field Experiments. *Ecol.*
749 *Monogr.* **54**, 187–211 (1984).
- 750 27. Ruxton, G. & Colegrave, N. *Experimental Design for the Life Sciences*. (OUP Oxford,
751 2011).
- 752 28. Stone, G. N., Nee, S. & Felsenstein, J. Controlling for non-independence in comparative
753 analysis of patterns across populations within species. *Philos. Trans. R. Soc. B Biol. Sci.*
754 **366**, 1410–1424 (2011).
- 755 29. Ives, A. R., Midford, P. E. & Garland, T., Jr. Within-Species Variation and Measurement
756 Error in Phylogenetic Comparative Methods. *Syst. Biol.* **56**, 252–270 (2007).
- 757 30. Bakkeren, E. *et al.* Cooperative virulence can emerge via horizontal gene transfer and is
758 stabilized by transmission. *Prep* (2021).
- 759 31. Ghoul, M., Andersen, S. B. & West, S. A. Sociomics: Using Omic Approaches to
760 Understand Social Evolution. *Trends Genet.* **33**, 408–419 (2017).
- 761 32. McInerney, J. O., McNally, A. & O’Connell, M. J. Why prokaryotes have pangenomes.
762 *Nat. Microbiol.* **2**, 17040 (2017).
- 763 33. Niehus, R., Mitri, S., Fletcher, A. G. & Foster, K. R. Migration and horizontal gene
764 transfer divide microbial genomes into multiple niches. *Nat. Commun.* **6**, (2015).
- 765 34. Cordero, O. X. *et al.* Ecological Populations of Bacteria Act as Socially Cohesive Units
766 of Antibiotic Production and Resistance. *Science* **337**, 1228–1231 (2012).
- 767 35. Rakoff-Nahoum, S., Coyne, M. J. & Comstock, L. E. An Ecological Network of
768 Polysaccharide Utilization among Human Intestinal Symbionts. *Curr. Biol.* **24**, 40–49
769 (2014).
- 770 36. Nocelli, N., Bogino, P. C., Banchio, E. & Giordano, W. Roles of Extracellular
771 Polysaccharides and Biofilm Formation in Heavy Metal Resistance of Rhizobia.
772 *Materials* **9**, 418 (2016).
- 773 37. Ciofu, O., Beveridge, T. J., Kadurugamuwa, J., Walther-Rasmussen, J. & Høiby, N.
774 Chromosomal β -lactamase is packaged into membrane vesicles and secreted from
775 *Pseudomonas aeruginosa*. *J. Antimicrob. Chemother.* **45**, 9–13 (2000).
- 776 38. Rodríguez-Beltrán, J., DelaFuente, J., León-Sampedro, R., MacLean, R. C. & San Millán,
777 Á. Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. *Nat. Rev.*
778 *Microbiol.* 1–13 (2021) doi:10.1038/s41579-020-00497-1.
- 779 39. Yu, N. Y. *et al.* PSORTb 3.0: improved protein subcellular localization prediction with
780 refined localization subcategories and predictive capabilities for all prokaryotes.
781 *Bioinformatics* **26**, 1608–1615 (2010).

- 782 40. Rankin, D. J., Rocha, E. P. C. & Brown, S. P. What traits are carried on mobile genetic
783 elements, and why? *Heredity* **106**, 1–10 (2011).
- 784 41. Hadfield, J. D. MCMC Methods for Multi-Response Generalized Linear Mixed Models:
785 The MCMCglmm R Package. *J. Stat. Softw.* **33**, 1–22 (2010).
- 786 42. Clutton-Brock, T. H. & Harvey, P. H. Primate ecology and social organization. *J. Zool.*
787 **183**, 1–39 (1977).
- 788 43. Crawley, M. J. *Statistics: An Introduction Using R*. (John Wiley & Sons, 2014).
- 789 44. Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering, reconstruction
790 and typing of plasmids from draft assemblies. *Microb. Genomics* **4**, (2018).
- 791 45. Robertson, J., Bessonov, K., Schonfeld, J. & Nash, J. H. E. Universal whole-sequence-
792 based plasmid typing and its utility to prediction of host range and epidemiological
793 surveillance. *Microb. Genomics* **6**, (2020).
- 794 46. Smillie, C., Garcillan-Barcia, M. P., Francia, M. V., Rocha, E. P. C. & de la Cruz, F.
795 Mobility of Plasmids. *Microbiol. Mol. Biol. Rev.* **74**, 434–452 (2010).
- 796 47. Mc Ginty, S. É. & Rankin, D. J. The evolution of conflict resolution between plasmids
797 and their bacterial hosts. *Evolution* **66**, 1662–1670 (2012).
- 798 48. Hamilton, W. D. Genetical evolution of social behaviour I & II. *J Theor Biol* **7**, 1–52
799 (1964).
- 800 49. work(s):, W. D. H. R. The Evolution of Altruistic Behavior. *Am. Nat.* **97**, 354–356
801 (1963).
- 802 50. Sheppard, R. J., Beddis, A. E. & Barraclough, T. G. The role of hosts, plasmids and
803 environment in determining plasmid transfer rates: A meta-analysis. *Plasmid* **108**,
804 102489 (2020).
- 805 51. Rodríguez-Beltrán, J. *et al.* Genetic dominance governs the evolution and spread of
806 mobile genetic elements in bacteria. *Proc. Natl. Acad. Sci.* **117**, 15755–15762 (2020).
- 807 52. Cornelis, G. R. *et al.* The Virulence Plasmid of Yersinia, an Antihost Genome. *Microbiol.*
808 *Mol. Biol. Rev.* **62**, 1315–1352 (1998).
- 809 53. Köstlbacher, S., Collingro, A., Halter, T., Domman, D. & Horn, M. Coevolving Plasmids
810 Drive Gene Flow and Genome Plasticity in Host-Associated Intracellular Bacteria. *Curr.*
811 *Biol.* **31**, 346-357.e3 (2021).
- 812 54. Gupta, A., Kapil, R., Dhakan, D. B. & Sharma, V. K. MP3: A Software Tool for the
813 Prediction of Pathogenic Proteins in Genomic and Metagenomic Data. *PLOS ONE* **9**,
814 e93907 (2014).

- 815 55. San Millan, A., Escudero, J. A., Gifford, D. R., Mazel, D. & MacLean, R. C. Multicopy
816 plasmids potentiate the evolution of antibiotic resistance in bacteria. *Nat. Ecol. Evol.* **1**,
817 0010 (2017).
- 818 56. Carrier, T., Jones, K. L. & Keasling, J. D. mRNA stability and plasmid copy number
819 effects on gene expression from an inducible promoter system. *Biotechnol. Bioeng.* **59**,
820 666–672 (1998).
- 821 57. Rodríguez-Beltrán, J. *et al.* Multicopy plasmids allow bacteria to escape from fitness
822 trade-offs during evolutionary innovation. *Nat. Ecol. Evol.* **2**, 873–881 (2018).
- 823 58. Dietel, A.-K., Kaltenpoth, M. & Kost, C. Convergent Evolution in Intracellular Elements:
824 Plasmids as Model Endosymbionts. *Trends Microbiol.* **26**, 755–768 (2018).
- 825 59. Rocha, E. P. C. & Danchin, A. Base composition bias might result from competition for
826 metabolic resources. *Trends Genet.* **18**, 291–294 (2002).
- 827 60. Garcia-Garcera, M., Touchon, M., Brisse, S. & Rocha, E. P. C. Metagenomic assessment
828 of the interplay between the environment and the genetic diversification of
829 *Acinetobacter*. *Environ. Microbiol.* **19**, 5010–5024 (2017).
- 830 61. Kümmerli, R., Schiessl, K. T., Waldvogel, T., McNeill, K. & Ackermann, M. Habitat
831 structure and the evolution of diffusible siderophores in bacteria. *Ecol. Lett.* **17**, 1536–
832 1544 (2014).
- 833 62. Canchaya, C., Fournous, G., Chibani-Chennoufi, S., Dillmann, M. L. & Brüssow, H.
834 Phage as agents of lateral gene transfer. *Curr. Opin. Microbiol.* **6**, 417–424 (2003).
- 835 63. Burrus, V. & Waldor, M. K. Shaping bacterial genomes with integrative and conjugative
836 elements. *Res. Microbiol.* **155**, 376–386 (2004).
- 837 64. O’Brien, F. G. *et al.* Origin-of-transfer sequences facilitate mobilisation of non-
838 conjugative antimicrobial-resistance plasmids in *Staphylococcus aureus*. *Nucleic Acids*
839 *Res.* **43**, 7971–7983 (2015).
- 840 65. Rodríguez-Rubio, L. *et al.* Extensive antimicrobial resistance mobilization via multicopy
841 plasmid encapsidation mediated by temperate phages. *J. Antimicrob. Chemother.* **75**,
842 3173–3180 (2020).
- 843 66. Ramsay, J. P. & Firth, N. Diverse mobilization strategies facilitate transfer of non-
844 conjugative mobile genetic elements. *Curr. Opin. Microbiol.* **38**, 1–9 (2017).
- 845 67. Jain, R., Rivera, M. C. & Lake, J. A. Horizontal gene transfer among genomes: the
846 complexity hypothesis. *Proc. Natl. Acad. Sci. U. S. A.* **96**, 3801–3806 (1999).

- 847 68. Cohen, O., Gophna, U. & Pupko, T. The complexity hypothesis revisited: connectivity
848 rather than function constitutes a barrier to horizontal gene transfer. *Mol. Biol. Evol.* **28**,
849 1481–1489 (2011).
- 850 69. Ding, W., Baumdicker, F. & Neher, R. A. panX: pan-genome analysis and exploration.
851 *Nucleic Acids Res.* **46**, e5 (2018).
- 852 70. Gardy, J. L. & Brinkman, F. S. L. Methods for predicting bacterial protein subcellular
853 localization. *Nat. Rev. Microbiol.* **4**, 741–751 (2006).
- 854 71. Ference, C. M. *et al.* Recent advances in the understanding of *Xanthomonas citri* ssp. *citri*
855 pathogenesis and citrus canker disease management. *Mol. Plant Pathol.* **19**, 1302–1318
856 (2018).
- 857 72. Morris, C. E., Lamichhane, J. R., Nikolić, I., Stanković, S. & Moury, B. The overlapping
858 continuum of host range among strains in the *Pseudomonas syringae* complex.
859 *Phytopathol. Res.* **1**, 4 (2019).
- 860 73. Hadfield, J. D. MCMCglmm Course Notes. Available at [cran.us.r-](http://cran.us.r-project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf)
861 [project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf](http://cran.us.r-project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf). (2019).
- 862 74. Nakagawa, S. & Schielzeth, H. A general and simple method for obtaining R² from
863 generalized linear mixed-effects models. *Methods Ecol. Evol.* **4**, 133–142 (2013).
- 864 75. Nakagawa, S., Johnson, P. C. D. & Schielzeth, H. The coefficient of determination R²
865 and intra-class correlation coefficient from generalized linear mixed-effects models
866 revisited and expanded. *J R Soc Interface* **11** (2017).
- 867 76. Hug, L. A. *et al.* A new view of the tree of life. *Nat. Microbiol.* **1**, 16048 (2016).
- 868 77. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and
869 evolutionary analyses in R. *Bioinforma. Oxf. Engl.* **35**, 526–528 (2019).
- 870 78. Revell, L. J. phytools: an R package for phylogenetic comparative biology (and other
871 things). *Methods Ecol. Evol.* **3**, 217–223 (2012).
- 872 79. Washburne, A. D. *et al.* Methods for phylogenetic analysis of microbiome data. *Nat.*
873 *Microbiol.* **3**, 652–661 (2018).
- 874 80. Som, A. Causes, consequences and solutions of phylogenetic incongruence. *Brief.*
875 *Bioinform.* **16**, 536–548 (2015).
- 876