

1 ***CRISPR-Cas is associated with fewer antibiotic resistance genes in bacterial***
2 ***pathogens***

3

4 Elizabeth Pursey^{1*}, Tatiana Dimitriu¹, Fernanda L. Paganelli², Edze R. Westra¹ and
5 Stineke van Houte¹

6

7 ¹Environment and Sustainability Institute, Biosciences, University of Exeter, Penryn,
8 Cornwall, United Kingdom

9 ²Department of Medical Microbiology, University Medical Center Utrecht, Utrecht,
10 The Netherlands

11

12 *Correspondence to: ellie.pursey@gmail.com

13

14 **Keywords:** CRISPR-Cas, horizontal gene transfer, mobile genetic elements,
15 antibiotic resistance, plasmids, integrative conjugative elements

16

17

18 **Abstract**

19

20 The acquisition of antibiotic resistance genes via horizontal gene transfer is a key
21 driver of the rise in multidrug resistance amongst bacterial pathogens. Bacterial
22 defence systems per definition restrict the influx of foreign genetic material, and may
23 therefore limit the acquisition of antibiotic resistance. CRISPR-Cas adaptive immune
24 systems are one of the most prevalent defences in bacteria, found in roughly half of
25 bacterial genomes, but it has remained unclear if and how much they contribute to
26 restricting the spread of antibiotic resistance. We analysed ~40,000 whole genomes
27 comprising the full RefSeq dataset for 11 species of clinically important genera of
28 human pathogens including *Enterococcus*, *Staphylococcus*, *Acinetobacter* and
29 *Pseudomonas*. We modelled the association between CRISPR-Cas and indicators of
30 horizontal gene transfer, and found that pathogens with a CRISPR-Cas system were
31 less likely to carry antibiotic resistance genes than those lacking this defence
32 system. Analysis of the mobile genetic elements targeted by CRISPR-Cas supports
33 a model where this host defence system blocks important vectors of antibiotic
34 resistance. These results suggest a potential “immunocompromised” state for
35 multidrug-resistant strains that may be exploited in tailored interventions that rely on
36 mobile genetic elements, such as phage or phagemids, to treat infections caused by
37 bacterial pathogens.

38

39 **Introduction**

40

41 The spread of antibiotic resistance (ABR) genes between bacterial strains and
42 species is of huge global importance; without access to working antibiotics, much of
43 modern medicine is threatened, including surgery, cancer treatment and neonatal
44 care [1]. Opportunistic pathogens in particular often have the ability to take up DNA
45 readily from the environment, through various mechanisms of horizontal gene
46 transfer (HGT) [2]. In some cases, strains are naturally competent and therefore able
47 to uptake extracellular DNA from their environment, as occurs in some clinical
48 isolates of *Acinetobacter baumannii* [3]. Mobile genetic elements (MGEs) also enter
49 the cell through transduction (transfer of DNA by bacteriophages), or conjugation, in
50 which mobile elements are physically transferred between cells through direct cell-
51 cell contact. Conjugative plasmids encode all necessary self-transfer machinery,

52 while mobilizable plasmids require it to be encoded by another element, but both
53 contribute hugely to the spread of ABR genes. For example, conjugative plasmids
54 are enriched for ABR genes and are able to transfer these resistance determinants
55 across species and genera [4]. In addition, they often harbour integrons,
56 characterised by the *int1* integrase gene, which enables them to capture ABR and
57 other genes in cassettes that can be spread by the conjugative plasmid [5].
58 Integrative conjugative elements (ICEs) are widespread in bacteria, and more
59 numerous than conjugative plasmids [6]. These MGEs can integrate into the host
60 chromosome and transmit vertically through cell division, but can also excise
61 themselves and be transferred horizontally through conjugation [7].

62

63 Host defences may block the acquisition of MGEs. One of the most prevalent
64 defences is CRISPR-Cas (clustered regularely interspaced short palindromic repeat
65 (CRISPR) loci and CRISPR-associated (*cas*) genes), an adaptive immune system
66 found in in ~30-40% of bacteria [8,9], in which short spacer sequences, derived from
67 incoming DNA, are incorporated at CRISPR loci in between repeat sequences.
68 Functioning as a form of immunological “memory”, these spacer sequences guide
69 their cognate Cas enzymes to cut and destroy any sequence matching this immune
70 record. This enables the cell to defend itself from predation by bacteriophages, and
71 block other sources of potentially costly incoming MGEs.

72

73 A number of studies have examined the extent to which CRISPR-Cas restricts HGT,
74 using genomic approaches. In contrast to virus-mediated immunity, which is often
75 circumvented through viral mutation, overcoming immunity to elements such as
76 plasmids often involves loss or inactivation of the CRISPR-Cas system itself [10].
77 This means evidence for its potential role in blocking these elements may be found
78 in the genomic and phylogenetic distributions of CRISPR-Cas systems in bacteria.

79

80 For example, some studies have examined whether CRISPR-Cas systems prevent
81 HGT by looking at genome length. Theoretically, genomes possessing CRISPR-Cas
82 may be shorter if they are blocking acquisition of foreign DNA elements. Evidence for
83 this has been found in *Pseudomonas aeruginosa* [11,12], a species with a large core
84 and accessory genome. More direct interactions between HGT and CRISPR-Cas
85 can be tested using genomic comparisons. However, the outcomes of these studies

86 have been ambiguous. For example, a study of 1399 bacterial and archaeal
87 genomes did not reveal a correlation between markers of recent HGT and spacer
88 count (used as a proxy for CRISPR-Cas activity) [13]. However, a more recent
89 analysis compared pairs of conspecific genomes with and without CRISPR-Cas as a
90 method to control for relatedness, and found fewer plasmids when CRISPR-Cas was
91 present within 29 genome pairs [14]. When comparing over 100,000 genomes from
92 multiple species, another study revealed a complex situation when looking for
93 correlations between particular ABR gene classes and CRISPR-Cas presence, in
94 which positive and negative relationships were found dependent on species and
95 ABR gene [15].

96

97 Research focused on within-species comparisons has often found more compelling
98 evidence for CRISPR-Cas as a barrier to HGT and the spread of ABR genes. In *P.*
99 *aeruginosa*, in paired within-ST (sequence type) genomes, fewer ICEs and
100 prophages were detected when CRISPR-Cas was present [12], whilst another study
101 found negative associations between the presence of type I-E and I-F systems and
102 particular ABR genes [11]. Studies in Enterococci also support the role of CRISPR-
103 Cas in blocking ABR acquisition, with fewer resistance genes detected in genomes
104 with CRISPR-Cas in a set of 48 strains of *E. faecalis* and 8 strains of *E. faecium* [16].
105 Further genomic associations between CRISPR-Cas presence and a lack of ABR
106 genes has been detected in *Klebsiella pneumoniae* [17,18], whilst there is evidence
107 in *A. baumannii* that CRISPR-Cas presence is correlated with an absence of
108 plasmids [19].

109

110 Here, we sought to expand on existing work using the complete RefSeq dataset for a
111 group of bacterial human pathogens from 11 species (n=39,511), with diverse
112 lifestyles, and CRISPR-Cas system types. We built upon the null hypothesis
113 significance testing used in previous studies, which typically used *t*-tests (or their
114 non-parametric alternatives) or correlation analyses. To do so, we applied linear and
115 generalised linear models (LMs, GLMs) with model selection based on Akaike
116 information criterion (AIC) values, which allows comparison of all candidate models
117 and selection of the model with the best fit to the dataset. We also used Bayesian
118 phylogenetic models to examine within-species relationships, to control for genetic
119 distance between genomes. Both of these model types are advantageous in that

120 they allow predictions to be made about the dataset based on values of different
121 variables.

122

123 We tested both between and within-species associations between ABR gene counts
124 and CRISPR-Cas presence. First, we assessed whether the probability that a
125 genome possesses a CRISPR-Cas system changes based on the presence of
126 plasmids, ICEs and integrons, as well as looking at whether genomes with CRISPR-
127 Cas were shorter. Subsequently, we built within-species models, controlling for
128 genetic distance, to detect associations between specific CRISPR-Cas types and
129 spacer counts, and the quantity of ABR genes accumulated in a genome. Finally, we
130 assessed whether a CRISPR-Cas system possessing spacers that target MGEs that
131 are vectors of ABR genes can effectively reduce the ABR count in the genome.

132

133 **Methods**

134

135 *Genomes*

136

137 RefSeq genomes for *P. aeruginosa*, *A. baumannii*, *Neisseria meningitidis*,
138 *Staphylococcus epidermidis*, *Streptococcus pyogenes*, *Francisella tularensis*,
139 *Mycobacterium tuberculosis*, *Neisseria gonorrhoeae*, *E. faecium* and *E. faecalis*
140 were retrieved from NCBI between December 2020 and January 2021 in nucleotide
141 fasta format using ncbi-genome-download v0.3.0 ([https://github.com/kbclin/ncbi-](https://github.com/kbclin/ncbi-genome-download)
142 [genome-download](https://github.com/kbclin/ncbi-genome-download)). The strains were selected to include a variety of CRISPR-Cas
143 system types, pathogenic lifestyles (i.e. obligate pathogens such as *F. tularensis*, *M.*
144 *tuberculosis* and *N. gonorrhoeae* as well as opportunistic pathogens like *P.*
145 *aeruginosa*, *S. aureus* and *S. pyogenes*), and species of significance in the spread of
146 ABR such as *A. baumannii* and *E. faecium*. The species used are summarised in
147 Table 1.

148

149 *Identification of CRISPR-Cas systems, ABR genes, plasmids, ICEs, and int1*

150

151 CRISPR-Cas systems were detected using CRISPRCasTyper [20], and orphan
152 CRISPR arrays were not included in any spacer analysis. CRISPR-Cas systems with
153 “Unknown” predictions were also removed from the dataset. Acquired ABR genes

154 and plasmid replicons were identified using abricate
155 (<https://github.com/tseemann/abricate>) with the NCBI AMRFinderPlus [21] and
156 PlasmidFinder [22] databases, both from 19/04/2019 and containing 5386 and 460
157 sequences, respectively. To detect ICEs and *intl1*, custom BLAST databases were
158 created based on the ICEberg database (retrieved January 2021; Liu et al. 2019),
159 and the NCBI gene entry for class 1 integron integrase *intl1* [NC_019069.1
160 (99852..100865, complement)], and used with abricate.

161

162 *Identification of spacer targets and MGEs carrying ABR genes*

163

164 Spacer sequences were searched against BLAST databases constructed from the
165 aforementioned ABR, ICE and prophage databases as well as a curated plasmid
166 database [24], and the complete RefSeq viral release (retrieved February 2021).
167 Only hits with 100% identity were used for downstream analyses. The total number
168 of spacers per genome was approximated by subtracting 1 from the number of
169 repeats for each array and summing these values, and any spacers with no hits in
170 any database were assigned as “unknown”.

171

172 MGEs carrying ABR genes were identified by running a blastn search with default
173 parameters, querying all sequences in search databases used for ICE, plasmid and
174 virus detection against a custom BLAST database built from the aforementioned
175 NCBI AMRFinderPlus database. Subsequently, a percent identity cutoff of 90% was
176 applied to hits. The presence of ABR genes on MGEs was classified in a binomial
177 way, in which ≥ 1 ABR gene was counted as the presence of ABR on the element. In
178 cases where multiple MGEs were targeted by a single spacer, the spacer was
179 classified as targeting an MGE with ABR if at least one of these MGEs carried at
180 least one ABR gene.

181

182 *Data analysis and statistical modelling*

183

184 The University of Exeter's Advanced Research Computing facilities were used to
185 carry out this work. Data analyses were conducted in R v4.0.2 using the tidyverse
186 suite of packages [25]. Unless otherwise specified, linear or generalized linear
187 models were fitted using the *lm* or *glm* functions in base R. Linear models allow the

188 inclusion of multiple predictors in one model, termed fixed effects, to estimate the
189 relationship between these predictors and a response variable. These models
190 generate a formula that calculates the mean of a response variable based on the
191 value of each fixed effect (predictor). Regression coefficients are fitted, which give
192 the slope and the intercept (value of the response variable when the predictor = 0)
193 for each fixed effect. For categorical predictors, a different intercept is fitted for each
194 level of the predictor (e.g. each species in a multispecies model). To fit different
195 slopes for levels of categorical predictor, interaction terms must be fitted between
196 them and a continuous predictor. For example, an interaction between species and
197 the count of ABR genes would allow varying slopes and intercepts for each individual
198 species.

199

200 For these analyses, a maximal model was generated and all possible candidate
201 models were compared using the AIC method using *dredge* from the MuMIn
202 package [26]. AIC values assess the fit of a model by looking at the likelihood of a
203 model given the data, penalising for increased number of parameters (as increased
204 complexity of the model increases parameter uncertainty). The model with the lowest
205 AIC value was selected, and no alternative models were within 2 AICs of the best
206 fitting model.

207

208 Prediction data frames with 95% confidence intervals were generated using
209 *ggpredict* from the package *ggeffects* [27], model dispersion was tested and scaled
210 residuals were examined using DHARMA residual diagnostics [28], and the final
211 predictions were visualised with *ggplot2*, *ghibli* [29] and *cowplot*
212 (<https://wilkelab.org/cowplot/index.html>). Prediction plots were produced only for
213 biologically realistic combinations of fixed effects and response, i.e. those where
214 sufficient data was present in the data set.

215

216 Associations between genome length and CRISPR-Cas presence/absence were
217 tested using a linear model with genome length as the response variable and
218 CRISPR-Cas presence/absence and species as fixed effects, with CRISPR-Cas
219 presence/absence \times species as an interaction term.

220

221 The effect of various MGEs on the distribution of CRISPR-Cas was modelled using a
222 binomial GLM with CRISPR-Cas presence/absence as the response variable and
223 counts of ICEs, plasmid replicons, ABR genes and *intl1* copies, as well as species,
224 as fixed effects. In addition, we fitted an interaction between species and every other
225 fixed effect. Prediction data frames were created up to the maximum number of each
226 MGE detected in each species (unless this value was an outlier, in which case a
227 more biologically informative upper limit was set).

228

229 The relationship between spacers targeting MGEs with and without ABR and the
230 count of ABR genes was modelled using a Bayesian Poisson GLM in MCMCglmm.
231 The default weakly informative priors were used, and models were run for 20,000
232 iterations with a burn in period of 10,000 iterations and a thinning interval of 5. Two
233 models were run, with the number of ABR genes as the response variable, species
234 as a fixed effect, and either the counts of spacers targeting MGEs with ABR or
235 targeting MGEs without ABR as an additional fixed effect. Prediction data frames
236 were created up to the maximum number of spacers detected targeting MGEs with
237 and without ABR.

238

239 *Construction of trees and genetic distance-controlled single-species models*

240

241 The package mashtree [30] with mash v2.0.0 [31] was used to create neighbour-
242 joining trees for each species based on whole genomes. Genomes were sketched
243 using mash with default parameters prior to using mashtree. Mash distance was
244 calculated for all isolates relative to the first genome entry for the group, in order to
245 remove outliers with a distance > 0.1, which are likely to have had their species
246 misclassified. Trees were subsequently rooted to outgroups and ultrametric trees
247 were produced with the ape package in R [32] using the chronos function with
248 smoothing parameters (λ) of 1 and 10. The “correlated”, “discrete” and “relaxed”
249 models were used for each value of λ , which vary in the extent to which they permit
250 adjacent parts of the tree to evolve at different rates. The tree with the highest log-
251 likelihood for each species was used. Due to the large number of genomes and low
252 prevalence of CRISPR-Cas, we did not perform this analysis for *S. aureus*.

253

254 The MCMCglmm package [33] was used to run Bayesian GLMs incorporating trees.
255 The default weakly informative priors were used, and all models were run for 30,000
256 iterations with a burn in period of 20,000 iterations and a thinning interval of 5.
257 Poisson GLMs were run for each species, with count of ABR genes as the response
258 variable and CRISPR-Cas type or spacer count (in which case genomes with no
259 CRISPR-Cas were removed) as a fixed effect. In cases where more than one
260 CRISPR-Cas type was present, CRISPR-Cas type \times total spacers was fit as an
261 interaction. The small sample size for *S. epidermidis* with CRISPR-Cas systems ($n =$
262 54), meant we chose not to model the association between spacer count and ABR
263 genes for this species. Prediction data frames for spacer counts were created up to
264 the maximum count of spacers detected for that species, with the exception of *S.*
265 *pyogenes*, for which the top value was an outlier at 33 so the next highest value of
266 17 was taken.

267

268 *Code availability and reproducibility*

269

270 The snakemake pipeline used to calculate genome lengths and detect CRISPR-Cas
271 systems, ABR genes, *int1*, ICEs and plasmid replicons, as well as R scripts used for
272 producing trees and statistical analysis for this work are available at
273 <https://github.com/elliepurple/crispr-pathogens>. Snakemake v5.18.1 [34] and
274 Python v3.8.3 with Biopython v1.78 [35] were used for this work.

275

276 **Results**

277

278 *Distribution of CRISPR-Cas systems, ABR and mobile elements across species*

279

280 Genome analyses of ~40,000 genomes of *Pseudomonas aeruginosa*, *Acinetobacter*
281 *baumannii*, *Neisseria meningitidis*, *Staphylococcus aureus*, *Staphylococcus*
282 *epidermidis*, *Streptococcus pyogenes*, *Francisella tularensis*, *Mycobacterium*
283 *tuberculosis*, *Neisseria gonorrhoeae*, *Enterococcus faecium* and *Enterococcus*
284 *faecalis* revealed a large variety of CRISPR-Cas system types, ABR genes, ICEs
285 and plasmids (Supp. Fig. 1). CRISPR-Cas systems were identified in every species
286 except *N. gonorrhoeae*, which was therefore excluded from subsequent analyses.

287 Overall, type I-C, I-E, I-F, II-A, II-B, II-C, III-A, IV and V-A systems were represented
288 across the dataset and ABR genes were detected in every species. However, *F.*
289 *tularensis* was positive only for its native β -lactamase (FTU-1) [36], *M. tuberculosis*
290 isolates were almost only ever positive for its intrinsic *blaA*, *aac(2')-Ic* and *erm(37)*
291 genes, and *N. meningitidis* had a maximum of 1 ABR gene, so these species were
292 not considered in downstream analyses. The proportion of genomes with CRISPR-
293 Cas varied across species from around 0.75 for *M. tuberculosis* to around 0.005 for
294 *S. aureus* (Supp. Fig. 2). The number of repeats per genome also varied
295 considerably between system type and species (Supp. Fig. 3), reaching as high as
296 ~200 in I-F systems in *A. baumannii*, and ~150 in III-A systems of *M. tuberculosis*.
297 However, for all other types this number was typically below 50.

298

299 Plasmids were detected in >80% of genomes for species in the *Staphylococcus* and
300 *Enterococcus* genera, and were also prevalent in *S. pyogenes* (~25%), *P.*
301 *aeruginosa* (10%), and *A. baumannii* (~3%) (Supp. Fig.1). As genomes were
302 assembled to varying levels within the dataset, it is assumed some plasmids will
303 have been missed using this method. ICEs were also common in these species,
304 except for *A. baumannii*, ranging from 93% of *E. faecalis* genomes having at least
305 one ICE to 15% of *E. faecium* genomes. Finally, *intI1* was detected in ~20% of *P.*
306 *aeruginosa* and *A. baumannii* genomes.

307

308 *No clear association between genome size and CRISPR-Cas presence/absence*

309

310 Consistent with previous work, *P. aeruginosa* genomes with CRISPR-Cas were on
311 average smaller than those without the system. This was in fact the largest
312 difference we saw, with predicted lengths of $6,579,416 \pm 2885$ bp and $6,692,831 \pm$
313 2601 bp for genomes with and without CRISPR-Cas, respectively (Supp. Fig. 4).
314 However, the opposite was true in other species such as *A. baumannii* and *S.*
315 *aureus*. For example, an *A. baumannii* genome with CRISPR-Cas was predicted to
316 be $4,028,659 \pm 5523$ bp long, whilst one without had a predicted length of $3,976,319$
317 ± 2181 bp, a difference of ~50,000 bp overall. Typically, differences between
318 CRISPR-Cas positive and negative genomes were < 55,000 bp.

319

320 *CRISPR-Cas presence is associated with fewer acquired ABR genes and mobile*
321 *elements, but more ICEs*

322

323 We modelled the association between CRISPR-Cas presence or absence and the
324 counts of ABR genes, as well as their vectors. Predictions are presented for
325 elements which are found in each species in Figure 1 and all model estimates are
326 presented in Supplementary Table 4. Across species, every additional ABR gene led
327 to a 0.08 reduction in the probability of having a CRISPR-Cas system ($p = 9.2E-18$).
328 The direction and strength of this trend varied by species, with positive associations
329 between probability of CRISPR-Cas presence and ABR gene count in *S. epidermidis*
330 and *P. aeruginosa* (Fig. 1A). CRISPR-Cas presence appeared to have little
331 association with plasmid replicon count across species (Fig. 1B), whilst for *intl1*, a
332 negative association was present in *P. aeruginosa* but not *A. baumannii*. For
333 genomes possessing ICEs, positive associations were detected between ICE count
334 and ABR gene count in *P. aeruginosa* and *S. pyogenes*.

335

336 *CRISPR-Cas presence is associated with fewer ABR genes across system types in*
337 *individual species, controlling for genetic distance*

338

339 To further explore the link between ABR genes and CRISPR-Cas, we made within-
340 species models controlling for genetic distance, to look at the effect of CRISPR-Cas
341 type and spacer count (Fig. 2). The overall trends detected are largely similar to
342 those presented for across species comparisons in Fig. 1, in that there is generally a
343 negative association between the presence of a CRISPR-Cas system and ABR
344 counts. Across species, lacking a CRISPR-Cas system was associated with a higher
345 predicted count of ABR genes (Fig. 2a). In some cases the size of this effect is
346 striking; for example, an *E. faecium* genome lacking CRISPR-Cas is predicted to
347 have 12 ABR genes compared to only 5 ABR genes for one possessing a type II-C
348 system. However, in most cases this effect is more modest, with a reduction of 1 or 2
349 ABR genes when CRISPR-Cas is present. There are, however, two notable
350 exceptions to this trend. Firstly, *Streptococcus pyogenes* type II-A systems were
351 associated with a small but nonsignificant increase (0.44, 95% CI = 0.35 - 0.56) in
352 the number of predicted ABR genes relative to genomes without CRISPR-Cas (0.33,

353 95% CI = 0.28 - 0.37). Secondly, both type I-C and IV systems are consistently
354 associated with an increased count of ABR genes.

355

356 Next, we modelled the number of ABR genes in response to the count of spacers
357 (for those genomes that possess CRISPR-Cas; Fig. 2b). Here we did not find strong
358 associations for most CRISPR-Cas system types. Exceptions were the type I-F
359 system of *A. baumannii*, which showed a dramatic reduction in ABR gene count with
360 increasing spacer count, and the type I-C system of *P. aeruginosa*, which
361 interestingly showed the opposite trend. We did not model this association for type
362 IV spacers due to the limited number of genomes (n = 20) in this category.

363

364 To expand on this dataset, we looked at the targets of all spacers by searching them
365 against virus, ICE and plasmid databases (Supp. Fig. 5). As expected, the majority
366 of the ~285,000 spacers detected (90%) did not have a match to any of these
367 elements. Overall, the majority of known spacer targets are viruses (n = 24,402).
368 However, we also detect many spacers targeting ICEs (n = 3,002) and plasmids (n =
369 1,416).

370

371 *Presence of spacers targeting ABR-carrying MGEs is negatively associated with*
372 *ABR genes*

373

374 As we detected many spacers matching potential vectors of ABR (n = 1,948), we
375 wanted to determine whether there was an association between the presence of
376 these spacers and an absence of ABR genes. We modelled the number of ABR
377 genes in response to the number of spacers targeting MGEs that were positive or
378 negative for ABR genes. The number of predicted ABR genes per genome was
379 consistently lower when spacers targeted MGEs carrying ABR than when they
380 targeted those that did not carry ABR, a trend that was universal across species (Fig.
381 3).

382

383 **Discussion**

384

385 We found little evidence for genome length as a marker of CRISPR-Cas blocking
386 HGT in this collection of bacterial pathogens. Genome reduction is a process that is

387 common in pathogens relative to their less-pathogenic relatives [37]. Bacteria may
388 streamline their genomes in order to save energetic resources and remove genes
389 that have become deleterious, or reductions in genome length may happen as a
390 result of many other evolutionary processes that occur during the transition to
391 pathogenicity. For example, in *P. aeruginosa*, strains adapted to the cystic fibrosis
392 lung are frequently auxotrophic, as they are able to utilise amino acids produced by
393 the host [38]. The extent of mobile DNA in genomes also reduces during the
394 transition from facultative to obligate lifestyles, potentially due to reduced
395 opportunities for their spread in more intracellular niches [39]. Therefore, this
396 measure may be detecting potential adaptation to a pathogenic lifestyle or particular
397 ecological environments, as well as a concurrent change in HGT frequencies.

398

399 However, we found compelling evidence that CRISPR-Cas can block the transfer of
400 MGEs. We focused on the spread of ABR and its potential vectors due to their global
401 importance, and identified a negative association between the count of ABR genes
402 and the probability a genome will have a CRISPR-Cas system, which was
403 repeatable across most CRISPR-Cas system types and species. Further supporting
404 this hypothesis, we found genomes with spacers targeting vectors of ABR did indeed
405 have fewer ABR genes. We did not detect strong trends for plasmid replicons,
406 potentially because 1) genomes were assembled to varying levels and sequenced
407 using different technologies within our dataset and many plasmids will have been
408 missed and 2) the database we used was built based on plasmids found in
409 *Enterobacteriaceae*, and has been shown to fail to predict plasmids in many species
410 [40].

411

412 When looking at specific CRISPR-Cas types within species, we found that ABR
413 genes were lower in frequency if CRISPR-Cas was present in most cases. Two
414 exceptions were type I-C and IV systems, which were associated with more ABR
415 genes. These system types may be driving the positive association we see between
416 CRISPR-Cas presence and ABR genes in our multi-species model for *P.*
417 *aeruginosa*. This result is also interesting in light of the fact that both of these system
418 types are present on MGEs. Type I-C systems occur on conjugative elements in *P.*
419 *aeruginosa* and have been found to correlate positively with the presence of
420 particular ABR genes [11], whilst type IV CRISPR-Cas systems are usually found on

421 plasmids, where they can mediate inter-plasmid competition [41]. CRISPR-Cas
422 system types that are more mobile may co-occur with ABR genes found on mobile
423 elements more frequently. Interestingly, the species for which we detected a positive
424 association between ICEs and CRISPR-Cas presence in our multi-species model (*P.*
425 *aeruginosa* and *S. pyogenes*) were also the only species to possess type I-C
426 systems. We also saw positive associations between the count of ABR genes and
427 the number of type I-C spacers in *P. aeruginosa*, suggesting that this system may
428 even allow the genome to accumulate more ABR genes, and could have a role in
429 competition between MGEs rather than host genome defence. However, much
430 remains to be learned about these widespread but relatively understudied type I-C
431 systems to fully explain this trend [42].

432

433 We do not see strong negative associations between the total count of spacers and
434 the presence of ABR genes across most species and CRISPR-Cas system types. As
435 the majority of spacer targets that we could identify were viruses, which are rarely
436 vectors of ABR [43], the large number of anti-phage spacers and comparative rarity
437 of spacers targeting ABR vectors may obscure any trend we would otherwise see
438 here. Alternatively, the lack of clear association may be explained by the tendency
439 for arrays to include more inactive spacers as they increase in size, as the most
440 recently acquired spacers are most active [44]. Therefore, genomes with more
441 spacers and larger arrays may include more spacers with no activity in removing
442 potential vectors of ABR. An exception to this trend was *A. baumannii*, which had a
443 strong reduction in ABR gene count with increasing spacer count in its type I-F
444 system. Our spacer target analysis found that, unusually, more spacers in this
445 species targeted plasmids than viruses; therefore, there could be some
446 specialisation towards blocking plasmid acquisition in this system, a trend which has
447 also been suggested for the type I-F system in *E. coli* [45].

448

449 We looked for spacers that targeted vectors of ABR and found they were fairly
450 common. An increase in the count of spacers targeting vectors of ABR had a striking
451 effect on reducing the predicted count of ABR genes. Particularly compared with the
452 weaker association between the count of spacers targeting vectors without ABR and
453 ABR genes, this evidence strongly suggests CRISPR-Cas is an important barrier to
454 vectors of ABR, thereby blocking the acquisition of ABR genes themselves. Many

455 ICEs, phages and plasmids remain to be characterised, and once our databases of
456 these elements expand, this hypothesis will become more easily testable using
457 genomic data.

458

459 This study supports the hypothesis that CRISPR-Cas system loss occurs where ABR
460 acquisition is beneficial. In this case, the presence of CRISPR-Cas may be selected
461 against in populations that are exposed to antibiotics, where the survival benefits of
462 acquiring ABR genes outweigh the cost of phage predation or MGE maintenance.
463 This has important implications for novel treatments, as multidrug-resistant bacteria
464 may represent an immunocompromised population. In recent years, increasing
465 attention has been given to antibiotic alternatives such as phage therapy [46,47] and
466 CRISPR-Cas antimicrobials delivered using phage-like or plasmid elements [48,49].
467 Our results support the use of these interventions in cases where antibiotics do not
468 work, as they are likely to be most effective in multidrug-resistant bacteria.

469

470 Our results also raise the question of why such trends have not been detected in
471 previous work asking similar questions, particularly when looking across species.
472 The strength of our work is primarily the use of more powerful and appropriate
473 statistical tests, as well as controls for genetic distance. We selected the model with
474 the best fit to the data, accounting and controlling for the effects of multiple
475 predictors under one modelling framework as well as testing interactions between
476 predictors. For example, we could test how the probability of a genome having a
477 CRISPR-Cas system varies based on the number of ABR genes it has, whilst
478 controlling for the number of ICEs, integrons and plasmid replicons in the genome
479 overall. We could also control for interactions; for example, this allowed us to detect
480 both positive and negative relationships between CRISPR-Cas presence and
481 genome length for different species within the same model, avoiding some of the
482 pitfalls of multiple testing. Such approaches are routinely used in evolution and
483 ecology studies to account for the variability of datasets [50], and while suitable for
484 these types of genomic comparisons, they are less commonly used in this context.
485 Importantly, we also have access to a much larger dataset of genomes and
486 improved computational tools since the publication of some previous work, which
487 allows us greater statistical power to uncover associations between CRISPR-Cas
488 and ABR genes.

489

490 **Conclusion**

491

492 We found that looking both across and within-species, CRISPR-Cas system
493 presence was associated with a reduction in counts of ABR genes. In addition,
494 where spacers targeted MGEs that carry ABR, a concurrent reduction in ABR genes
495 was seen. These results have promising implications for the delivery of novel
496 technologies to combat ABR such as phage therapy and CRISPR-Cas
497 antimicrobials, which rely on bypassing bacterial immune systems to kill cells. In fact,
498 they may be ideal for this purpose, if the most multidrug-resistant strains are also the
499 most immunocompromised.

500

501 **Figure legends**

502

503 Figure 1 – Predicted probability of CRISPR-Cas presence (y-axis) from binomial
504 GLMs for A) ABR gene count and numbers of the potential ABR vectors B) plasmids,
505 C) *intl1* (integrons) and D) ICEs. Predictions are presented for elements which are
506 found in each species. Error bars show 95% confidence intervals.

507

508 Figure 2 – Prediction plots from Bayesian Poisson GLMs of ABR gene counts for A)
509 CRISPR-Cas type and B) count of spacers. Each column represents an individual
510 species. Error bars show 95% credible intervals.

511

512 Figure 3 – Prediction plot overlaying two Bayesian Poisson GLMs of ABR gene
513 counts according to counts of spacers (with known targets) targeting MGEs with and
514 without ABR, faceted by species. Error bars show 95% credible intervals.

515

516 **Table legend**

517

518 Table 1 – Summary of species included

519

520 **Acknowledgements**

521

522 We would like to thank MD Sharma for his invaluable help with accessing
523 computational resources needed for this work.

524

525 **Funding**

526

527 EP is supported by a PhD studentship equally funded by the European Research
528 Council (ERC-STG-2016-714478 – EVOIMMECH) and the College of Life and
529 Environmental Sciences, University of Exeter. SVH acknowledges support from the
530 Biotechnology and Biological Sciences Research Council (BB/S017674/1 and
531 BB/R010781/10). FLP acknowledges support from Utrecht Exposome Hub of Utrecht
532 Life Sciences (www.uu.nl/exposome), funded by the Executive Board of Utrecht
533 University.

534

535 **Competing Interests**

536

537 We have no competing interests.

538

539 **Table**

Species	Abbreviation used	RefSeq genomes (n)
<i>Pseudomonas aeruginosa</i>	PA	5360
<i>Acinetobacter baumannii</i>	AB	4864
<i>Neisseria meningitidis</i>	NM	1967
<i>Staphylococcus epidermidis</i>	SE	908
<i>Staphylococcus aureus</i>	SA	12148
<i>Streptococcus pyogenes</i>	SP	2106
<i>Francisella tularensis</i>	FT	675
<i>Mycobacterium tuberculosis</i>	MT	6688
<i>Neisseria gonorrhoeae</i>	NG	723
<i>Enterococcus faecium</i>	EFm	2247
<i>Enterococcus faecalis</i>	EFs	1825
	Total	39511

540

541 **Bibliography**

542

- 543 1. Laxminarayan R, Matsoso P, Pant S, Brower C, Røttingen J-A, Klugman K,
544 Davies S. 2016 Access to effective antimicrobials: a worldwide challenge. *The*
545 *Lancet* **387**, 168–175. (doi:10.1016/S0140-6736(15)00474-2)
- 546 2. Von Wintersdorff CJH, Penders J, Van Niekerk JM, Mills ND, Majumder S, Van
547 Alphen LB, Savelkoul PHM, Wolffs PFG. 2016 Dissemination of antimicrobial
548 resistance in microbial ecosystems through horizontal gene transfer. *Front.*
549 *Microbiol.* **7**, 1–10. (doi:10.3389/fmicb.2016.00173)
- 550 3. Traglia GM, Chua K, Centrón D, Tolmasky ME, Ramírez MS. 2014 Whole-
551 Genome Sequence Analysis of the Naturally Competent *Acinetobacter*
552 *baumannii* Clinical Isolate A118. *Genome Biol. Evol.* **6**, 2235–2239.
553 (doi:10.1093/gbe/evu176)
- 554 4. Che Y, Yang Y, Xu X, Břinda K, Polz MF, Hanage WP, Zhang T. 2021
555 Conjugative plasmids interact with insertion sequences to shape the horizontal
556 transfer of antimicrobial resistance genes. *Proc. Natl. Acad. Sci.* **118**,
557 e2008731118. (doi:10.1073/pnas.2008731118)
- 558 5. Vrancianu CO, Popa LI, Bleotu C, Chifiriuc MC. 2020 Targeting Plasmids to Limit
559 Acquisition and Transmission of Antimicrobial Resistance. *Front. Microbiol.* **11**,
560 761. (doi:10.3389/fmicb.2020.00761)
- 561 6. Guglielmini J, Quintais L, Garcillán-Barcia MP, de la Cruz F, Rocha EPC. 2011
562 The Repertoire of ICE in Prokaryotes Underscores the Unity, Diversity, and
563 Ubiquity of Conjugation. *PLoS Genet.* **7**, e1002222.
564 (doi:10.1371/journal.pgen.1002222)
- 565 7. Botelho J, Schulenburg H. 2021 The Role of Integrative and Conjugative
566 Elements in Antibiotic Resistance Evolution. *Trends Microbiol.* **29**, 8–18.
567 (doi:10.1016/j.tim.2020.05.011)
- 568 8. Hatoum-Aslan A, Marraffini LA. 2014 Impact of CRISPR immunity on the
569 emergence and virulence of bacterial pathogens. *Curr. Opin. Microbiol.* **17**, 82–
570 90. (doi:10.1016/j.mib.2013.12.001)
- 571 9. Bernheim A, Bikard D, Touchon M, Rocha EPC. 2019 Atypical organizations and
572 epistatic interactions of CRISPRs and cas clusters in genomes and their mobile
573 genetic elements. *Nucleic Acids Res.* , gkz1091. (doi:10.1093/nar/gkz1091)
- 574 10. Jiang W, Maniv I, Arain F, Wang Y, Levin BR, Marraffini LA. 2013 Dealing with
575 the Evolutionary Downside of CRISPR Immunity: Bacteria and Beneficial
576 Plasmids. *PLoS Genet.* **9**, e1003844. (doi:10.1371/journal.pgen.1003844)
- 577 11. van Belkum A *et al.* 2015 Phylogenetic Distribution of CRISPR-Cas Systems in
578 Antibiotic-Resistant *Pseudomonas aeruginosa*. **6**, 13.
- 579 12. Wheatley RM, MacLean RC. 2020 CRISPR-Cas systems restrict horizontal gene
580 transfer in *Pseudomonas aeruginosa*. *ISME J.* (doi:10.1038/s41396-020-00860-
581 3)

- 582 13. Gophna U, Kristensen DM, Wolf YI, Popa O, Drevet C, Koonin EV. 2015 No
583 evidence of inhibition of horizontal gene transfer by CRISPR–Cas on
584 evolutionary timescales. *ISME J.* **9**, 2021–2027. (doi:10.1038/ismej.2015.20)
- 585 14. O’Meara D, Nunney L. 2019 A phylogenetic test of the role of CRISPR-Cas in
586 limiting plasmid acquisition and prophage integration in bacteria. *Plasmid* **104**,
587 102418. (doi:10.1016/j.plasmid.2019.102418)
- 588 15. Shehreen S, Chyou T, Fineran PC, Brown CM. 2019 Genome-wide correlation
589 analysis suggests different roles of CRISPR-Cas systems in the acquisition of
590 antibiotic resistance genes in diverse species. *Philos. Trans. R. Soc. B Biol. Sci.*
591 **374**, 20180384. (doi:10.1098/rstb.2018.0384)
- 592 16. Palmer KL, Gilmore MS. 2010 Multidrug-Resistant Enterococci Lack CRISPR-
593 cas. *mBio* **1**, e00227-10. (doi:10.1128/mBio.00227-10)
- 594 17. Mackow NA, Shen J, Adnan M, Khan AS, Fries BC, Diago-Navarro E. 2019
595 CRISPR-Cas influences the acquisition of antibiotic resistance in *Klebsiella*
596 *pneumoniae*. *PLOS ONE* **14**, e0225131. (doi:10.1371/journal.pone.0225131)
- 597 18. Li H-Y, Kao C-Y, Lin W-H, Zheng P-X, Yan J-J, Wang M-C, Teng C-H, Tseng C-
598 C, Wu J-J. 2018 Characterization of CRISPR-Cas Systems in Clinical *Klebsiella*
599 *pneumoniae* Isolates Uncovers Its Potential Association With Antibiotic
600 Susceptibility. *Front. Microbiol.* **9**. (doi:10.3389/fmicb.2018.01595)
- 601 19. Mangas EL, Rubio A, Álvarez-Marín R, Labrador-Herrera G, Pachón J, Pachón-
602 Ibáñez ME, Divina F, Pérez-Pulido AJ. 2019 Pangenome of *Acinetobacter*
603 *baumannii* uncovers two groups of genomes, one of them with genes involved in
604 CRISPR/Cas defence systems associated with the absence of plasmids and
605 exclusive genes for biofilm formation. *Microb. Genomics* **5**.
606 (doi:10.1099/mgen.0.000309)
- 607 20. Russel J, Pinilla-Redondo R, Mayo-Muñoz D, Shah SA, Sørensen SJ. 2020
608 CRISPRCasTyper: Automated Identification, Annotation, and Classification of
609 CRISPR-Cas Loci. *CRISPR J.* **3**, 462–469. (doi:10.1089/crispr.2020.0059)
- 610 21. Feldgarden M *et al.* 2019 Validating the AMRFinder Tool and Resistance Gene
611 Database by Using Antimicrobial Resistance Genotype-Phenotype Correlations
612 in a Collection of Isolates. *Antimicrob. Agents Chemother.* **63**.
613 (doi:10.1128/AAC.00483-19)
- 614 22. Carattoli A, Zankari E, García-Fernández A, Voldby Larsen M, Lund O, Villa L,
615 Møller Aarestrup F, Hasman H. 2014 *In Silico* Detection and Typing of Plasmids
616 using PlasmidFinder and Plasmid Multilocus Sequence Typing. *Antimicrob.*
617 *Agents Chemother.* **58**, 3895–3903. (doi:10.1128/AAC.02412-14)
- 618 23. Liu M, Li X, Xie Y, Bi D, Sun J, Li J, Tai C, Deng Z, Ou H-Y. 2019 ICEberg 2.0:
619 an updated database of bacterial integrative and conjugative elements. *Nucleic*
620 *Acids Res.* **47**, D660–D665. (doi:10.1093/nar/gky1123)

- 621 24. Brooks L, Kaze M, Siström M. 2019 A Curated, Comprehensive Database of
622 Plasmid Sequences. *Microbiol. Resour. Announc.* **8**, e01325-18, e01325-18.
623 (doi:10.1128/MRA.01325-18)
- 624 25. Wickham H *et al.* 2019 Welcome to the Tidyverse. *J. Open Source Softw.* **4**,
625 1686. (doi:10.21105/joss.01686)
- 626 26. Barton K. 2009 MuMIn: multi-model inference. *Http-Forg.-Proj.*
- 627 27. Lüdtke D. 2018 ggeffects: Tidy Data Frames of Marginal Effects from
628 Regression Models. *J. Open Source Softw.* **3**, 772. (doi:10.21105/joss.00772)
- 629 28. Hartig F. 2020 *DHARMA: Residual Diagnostics for Hierarchical (Multi-Level /*
630 *Mixed) Regression Models*. See <http://florianhartig.github.io/DHARMA/>.
- 631 29. Henderson E. 2020 *ghibli: Studio Ghibli Colour Palettes*. See [https://CRAN.R-](https://CRAN.R-project.org/package=ghibli)
632 [project.org/package=ghibli](https://CRAN.R-project.org/package=ghibli).
- 633 30. Katz LS, Griswold T, Morrison SS, Caravas JA, Zhang S, Bakker HC den, Deng
634 X, Carleton HA. 2019 Mashtree: a rapid comparison of whole genome sequence
635 files. *J. Open Source Softw.* **4**, 1762. (doi:10.21105/joss.01762)
- 636 31. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S,
637 Phillippy AM. 2016 Mash: fast genome and metagenome distance estimation
638 using MinHash. *Genome Biol.* **17**, 132. (doi:10.1186/s13059-016-0997-x)
- 639 32. Paradis E, Claude J, Strimmer K. 2004 APE: Analyses of Phylogenetics and
640 Evolution in R language. *Bioinformatics* **20**, 289–290.
641 (doi:10.1093/bioinformatics/btg412)
- 642 33. Hadfield JD. 2010 MCMC Methods for Multi-Response Generalized Linear Mixed
643 Models: The MCMCglmm R Package. *J. Stat. Softw.* **33**.
644 (doi:10.18637/jss.v033.i02)
- 645 34. Köster J, Rahmann S. 2012 Snakemake—a scalable bioinformatics workflow
646 engine. *Bioinformatics* **28**, 2520–2522. (doi:10.1093/bioinformatics/bts480)
- 647 35. Cock PJA *et al.* 2009 Biopython: freely available Python tools for computational
648 molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423.
649 (doi:10.1093/bioinformatics/btp163)
- 650 36. Antunes NT, Frase H, Toth M, Vakulenko SB. 2012 The Class A β -Lactamase
651 FTU-1 Is Native to *Francisella tularensis*. *Antimicrob. Agents Chemother.* **56**,
652 666–671. (doi:10.1128/AAC.05305-11)
- 653 37. Moran NA. 2002 Microbial Minimalism: Genome Reduction in Bacterial
654 Pathogens. *Cell* **108**, 583–586. (doi:10.1016/S0092-8674(02)00665-7)
- 655 38. Weinert LA, Welch JJ. 2017 Why Might Bacterial Pathogens Have Small
656 Genomes? *Trends Ecol. Evol.* **32**, 936–947. (doi:10.1016/j.tree.2017.09.006)

- 657 39. Newton ILG, Bordenstein SR. 2011 Correlations Between Bacterial Ecology and
658 Mobile DNA. *Curr. Microbiol.* **62**, 198–208. (doi:10.1007/s00284-010-9693-3)
- 659 40. Arredondo-Alonso S, Willems RJ, van Schaik W, Schürch AC. 2017 On the
660 (im)possibility of reconstructing plasmids from whole-genome short-read
661 sequencing data. *Microb. Genomics* **3**. (doi:10.1099/mgen.0.000128)
- 662 41. Pinilla-Redondo R, Mayo-Muñoz D, Russel J, Garrett RA, Randau L, Sørensen
663 SJ, Shah SA. 2020 Type IV CRISPR–Cas systems are highly diverse and
664 involved in competition between plasmids. *Nucleic Acids Res.* **48**, 2000–2012.
665 (doi:10.1093/nar/gkz1197)
- 666 42. León LM, Park AE, Borges AL, Zhang JY, Bondy-Denomy J. 2021 Mobile
667 element warfare via CRISPR and anti-CRISPR in *Pseudomonas aeruginosa*.
668 *Nucleic Acids Res.* **49**, 2114–2125. (doi:10.1093/nar/gkab006)
- 669 43. Enault F, Briet A, Bouteille L, Roux S, Sullivan MB, Petit M-A. 2017 Phages
670 rarely encode antibiotic resistance genes: a cautionary tale for virome analyses.
671 *ISME J.* **11**, 237–247. (doi:10.1038/ismej.2016.90)
- 672 44. McGinn J, Marraffini LA. 2016 CRISPR-Cas Systems Optimize Their Immune
673 Response by Specifying the Site of Spacer Integration. *Mol. Cell* **64**, 616–623.
674 (doi:10.1016/j.molcel.2016.08.038)
- 675 45. Aydin S, Personne Y, Newire E, Laverick R, Russell O, Roberts AP, Enne VI.
676 2017 Presence of Type I-F CRISPR/Cas systems is associated with antimicrobial
677 susceptibility in *Escherichia coli*. *J. Antimicrob. Chemother.* **72**, 2213–2218.
678 (doi:10.1093/jac/dkx137)
- 679 46. Lin DM, Koskella B, Lin HC. 2017 Phage therapy: An alternative to antibiotics in
680 the age of multi-drug resistance. *World J. Gastrointest. Pharmacol. Ther.* **8**, 162–
681 173. (doi:10.4292/wjgpt.v8.i3.162)
- 682 47. Pires DP, Costa AR, Pinto G, Meneses L, Azeredo J. 2020 Current challenges
683 and future opportunities of phage therapy. *FEMS Microbiol. Rev.* **44**, 684–700.
684 (doi:10.1093/femsre/fuaa017)
- 685 48. Bikard D, Barrangou R. 2017 Using CRISPR-Cas systems as antimicrobials.
686 *Curr. Opin. Microbiol.* **37**, 155–160. (doi:10.1016/j.mib.2017.08.005)
- 687 49. Pursey E, Sünderhauf D, Gaze W, Westra ER, van Houte S. 2018 CRISPR-Cas
688 antimicrobials: Challenges and future prospects. *PLoS Pathog.* **14**, e1006990.
- 689 50. Bolker BM, Brooks ME, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, White
690 J-SS. 2009 Generalized linear mixed models: a practical guide for ecology and
691 evolution. *Trends Ecol. Evol.* **24**, 127–135. (doi:10.1016/j.tree.2008.10.008)





