

1 **Timescale and genetic linkage explain the variable impact of defense systems on horizontal gene**
2 **transfer**

3 Yang Liu¹, João Botelho¹, Jaime Iranzo^{2,3}

4 ¹ Centro de Biotecnología y Genómica de Plantas, Universidad Politécnica de Madrid (UPM) -
5 Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria (INIA-CSIC), Madrid, Spain.

6 ² Centro de Astrobiología (CAB), CSIC-INTA, Madrid, Spain.

7 ³ Institute for Biocomputation and Physics of Complex Systems (BIFI), University of Zaragoza,
8 Zaragoza, Spain.

9

10 **Running title:** Defense systems and horizontal gene transfer

11

12 **Keywords:** defense system, horizontal gene transfer, mobile genetic element, CRISPR-Cas, phage,
13 plasmid, comparative genomics.

14 **Abstract**

15 Prokaryotes have evolved a wide repertoire of defense systems to prevent invasion by mobile
16 genetic elements (MGE). However, because MGE are vehicles for the exchange of beneficial
17 accessory genes, defense systems could consequently impede rapid adaptation in microbial
18 populations. Here, we study how defense systems impact horizontal gene transfer (HGT) in the short
19 and long terms. By combining comparative genomics and phylogeny-aware statistical methods, we
20 quantified the association between the presence of 7 widespread defense systems and the
21 abundance of MGE in the genomes of 196 bacterial and 1 archaeal species. We also calculated the
22 differences in the rates of gene gain and loss between lineages that possess and lack each defense
23 system. Our results show that the impact of defense systems on HGT is highly taxon- and system-
24 dependent, and in most cases not statistically significant. Timescale analysis reveals that defense
25 systems must persist in a lineage for a relatively long time to exert an appreciable negative impact on
26 HGT. In contrast, for shorter evolutionary timescales, frequent co-acquisition of MGE and defense
27 systems results in a net positive association of the latter with HGT. Given the high turnover rates
28 experienced by defense systems, we propose that the inhibitory effect of most defense systems on
29 HGT is masked by their strong linkage with MGE. These findings help explain the contradictory
30 conclusions of previous research by pointing at mobility and within-host retention times as key
31 factors that determine the impact of defense systems on genome plasticity.

32 **Introduction**

33 Gene exchange plays a key role in the adaptation of microbes to changing environments, facilitating
34 the spread of antibiotic resistance, pathogenicity factors, metabolic genes, and other accessory
35 functions (Arnold et al. 2022). Over the last decade, there has been an increasing interest in assessing
36 the ecological and genetic factors that control horizontal gene transfer (HGT) and determine the
37 outcome of newly acquired genes in microbial populations (Soucy et al. 2015; Hall et al. 2020; Lee et

38 al. 2022). HGT is often mediated by mobile genetic elements (MGE), such as phages, integrative and
39 conjugative elements, and plasmids, against which bacteria have evolved an elaborate repertoire of
40 defense systems (Doron et al. 2018; Botelho et al. 2023; Georjon and Bernheim 2023; Mayo-Munoz
41 et al. 2023; Shaw et al. 2023). As a result, large-scale patterns of HGT are shaped by an interplay of
42 ecological and genetic variables that underlie cross-strain and cross-species differences in
43 susceptibility to MGE (Haudiquet et al. 2022).

44 Recent studies have highlighted the role of CRISPR-Cas as widespread adaptive immunity systems
45 that protect archaea and bacteria against viruses and other MGE (van Vliet et al. 2021; Watson et al.
46 2024). While in vitro experiments have provided supportive evidence for this function (Marraffini and
47 Sontheimer 2008; O'Hara et al. 2017; Watson et al. 2018), questions remain about the extent to
48 which CRISPR-Cas systems constrain gene exchange in nature (Gophna et al. 2015; O'Meara and
49 Nunney 2019; Shehreen et al. 2019; Westra and Levin 2020; Wheatley and MacLean 2021; Pursey et
50 al. 2022). More generally, there is a paucity of research on how other defense systems, such as
51 restriction-modification (RM), abortive infection (Abi), and an expansive repertoire of recently
52 identified gene systems including Gabija, CBASS (cyclic oligonucleotide-based antiphage signaling
53 system), DMS (DNA modification-based systems), and DRT (defense-associated reverse
54 transcriptases) affect HGT in prokaryotes (Tesson et al. 2022; Costa et al. 2024).

55 In contrast with the expectation that defense systems restrict HGT by interfering with the
56 propagation of MGE, several empirical and theoretical observations suggest that the relation
57 between defense systems and HGT might be more complex. First, the selection pressure to maintain
58 defense systems in a population (and consequently their prevalence) generally increases with the
59 exposure to MGE (Oliveira et al. 2016; Meaden et al. 2022). Second, defense systems are mobilized
60 by MGE, which could lead to a trivial positive association with HGT rates. In a less trivial manner,
61 whole defense systems or parts of them are often encoded by MGE (Makarova et al. 2011; Pinilla-

62 Redondo et al. 2022; Botelho 2023), promoting the retention of the latter for their beneficial side-
63 effects in cellular defense (Koonin et al. 2020; Rocha and Bikard 2022).

64 Here, we investigated the association between 7 widespread defense systems and HGT rates in 197
65 prokaryotic species. By combining high-quality genomic data, phylogenomic methods, and
66 phylogeny-aware statistical inference, our study aimed to shed light on the nuanced consequences of
67 the interplay between defense systems and MGE on bacterial evolution across timescales.

68

69 **Results**

70 Association between defense systems and MGE abundance is MGE- and taxon-dependent

71 We used species-wise phylogenetic generalized linear mixed models (PGLMM) to study the
72 association between the presence or absence of the 7 most prevalent defense systems in the dataset
73 (RM, DMS, Abi, CRISPR-Cas, Gabija, DRT, and CBASS), genome size (measured as the total number of
74 genes), and the number of MGE per genome (Fig. 1A-B, Tables S1-S4). Notably, the sign and
75 magnitude of the associations are strongly taxon-, system-, and MGE-dependent. Out of 197 species
76 included in the analysis, around 20-30 (depending on the defense system) displayed a statistically
77 significant ($p < 0.05$) positive association between the presence of the defense system and genome
78 size. In contrast, statistically significant negative associations were only observed in 5-15 species (Fig.
79 1A, top row). Differences in the number of genomes per species could bias the assessment of
80 statistical significance, leading to an overrepresentation of well-sampled taxa among those that
81 display significant associations. Moreover, the detection of associations involving extremely
82 abundant systems, such as RM and DMS (both with mean within-species prevalence close to 90%),
83 could be compromised by insufficient statistical power. To overcome these limitations, we
84 implemented a more permissive (but less biased) alternative criterion based on effect sizes to
85 determine the number of positive and negative associations (see Methods). Regardless of the

86 criterion, the presence of defense systems only correlates with genome size and MGE abundance in a
87 minority of species. Overall, positive associations outnumber negative associations (2:1 and 1.5:1
88 positive-vs-negative ratios for genome size and MGE abundance, according to the effect size
89 criterion; 3.5:1 and 2:1 with statistical significance), although the trend is less pronounced in CRISPR-
90 Cas systems (1.17:1 and 1:1 for genome size and MGE abundance; 1.75:1 and 1.5:1 with statistical
91 significance). The analysis also reveals differences regarding the association between defense
92 systems and distinct types of MGE (Fig. 1A, middle rows). In all defense systems, negative
93 associations with prophages are more frequent than negative associations with plasmids (1.2-1.9
94 times more frequent, according to the effect size criterion; 1.8-4.7 with the statistical significance
95 criterion). This is especially manifest in the case of CRISPR-Cas, whose presence correlates with a
96 reduction in the number of prophages in 68 species (14 statistically significant), and with higher
97 numbers of transposable elements and plasmids in 64 and 58 species, respectively (18 and 11
98 statistically significant). Significant associations, when detected, affect a sizeable fraction of the
99 accessory genome, with 20-40% differences in MGE content between genomes that do and do not
100 harbor the defense system (Fig. 1B).

101 A more detailed analysis at the level of functional categories reveals that the presence of defense
102 systems is most often associated with changes in the number of genes from COG categories X
103 (mobilome), L (replication, recombination and repair), U (intracellular trafficking and secretion), K
104 (transcription), D (cell cycle control, cell division and chromosome partitioning), and V (defense) (Fig.
105 1A and S1). Genes from these functional categories are typically present in MGE, suggesting that
106 correlations (both positive and negative) between defense systems and genome size are primarily
107 due to differences in the abundance of MGE. We confirmed that by masking genomic regions that
108 correspond to known MGE and rerunning the statistical analysis. As expected, 60% of the significant
109 associations disappeared after masking MGE (Fig S2). Although some significant associations
110 persisted, a closer inspection revealed that those were often related to genes from degenerated

111 prophages and other MGE (such as phage satellites) that had not been originally identified as such
112 and remained unmasked.

113 The sign of the association between defense systems and MGE does not follow a clear taxonomic
114 trend (Fig. 2, S3 and S4), with opposite signs sometimes found in closely related species (see, for
115 example, the differences for CRISPR-Cas in *Phocaeicola vulgatus* and *P. dorei*). Moreover, the same
116 species often display opposite trends for different defense systems. For example, in *Pseudomonas*
117 *aeruginosa*, genomes with CRISPR-Cas contain fewer MGE, whereas genomes with CBASS, Gabija,
118 and RM systems are significantly enriched in MGE. Negative associations between CRISPR-Cas and
119 MGE are more abundant in the genus *Acinetobacter* (5 out of 8 species), the phylum *Bacteroidota*
120 (negative association in 8 species, positive association in a single species) and the class *Clostridia*, the
121 latter especially affecting prophages (negative association in 11 species, positive association in 2
122 species; Fig. S4). The genus *Acinetobacter* is also enriched in negative associations involving Abi (5
123 species) and CBASS (4 species). Furthermore, three almost non-overlapping groups of streptococci
124 display negative associations for different defense systems. These encompass *S. pyogenes*, *S.*
125 *gordonii*, *S. anginosus*, *S. mutans*, and *S. salivarius* in the case of CRISPR-Cas; *S. anginosus*, *S. oralis*, *S.*
126 *intermedius*, and *S. suis* in the case of RM and DMS; and *S. pyogenes*, *S. dysgalactiae*, *S. equi*, and *S.*
127 *uberis* in the case of Gabija.

128 The only archaeal species in the study (*Methanococcus maripaludis*) did not show any remarkable
129 trend, other than a relatively strong (but not statistically significant) positive association between
130 prophages and DMS/RM, and between transposons and Gabija, and a moderate negative association
131 between prophages and Gabija.

132 Associations between defense systems and MGE arise from differences in the rate of gene
133 acquisition

134 To investigate if correlations between defense systems and MGE involve cross-strain differences in
135 genome plasticity, we built high-resolution strain trees and identified clades that contain defense

136 systems (DEF⁺). Then, we inferred the rates of gene gain and loss associated with different classes of
137 MGE in DEF⁺ clades and in their respective sister groups lacking the defense system (DEF⁻) using the
138 phylogenomic reconstruction tool GLOOME. We quantified the relative differences in gene gain and
139 loss by dividing the rates observed in DEF⁺ clades by those from sister DEF⁻ clades. Finally, we
140 compared the resulting DEF⁺/DEF⁻ gain and loss ratios between species in which MGE abundances
141 and defense systems are positively and negatively correlated. We found that DEF⁺ and DEF⁻ clades
142 often differ in their gain rates but not in their loss rates (Fig. 3A-B). Specifically, in species in which
143 defense systems are associated with increased MGE abundance, gene gain rates are typically higher
144 in clades that contain the defense system. The opposite (that is, a reduction of gene gain rates in
145 DEF⁺ clades) is observed in species in which defense systems are associated with lower MGE
146 abundance. Among the latter, the biggest reductions in plasmid acquisition occur in association with
147 RM, DRT, DMS and Abi, whereas significant drops in prophage gain are observed for DRT, CBASS and
148 CRISPR-Cas. Taken together, these results confirm that correlations between defense systems and
149 MGE arise from cross-strain differences in gene gain rather than loss, as expected given the role of
150 MGE as facilitators of HGT.

151 Timescales and linkage determine the sign of associations between defense systems and HGT

152 Because prokaryotic defense systems are often located within or adjacent to MGE, we hypothesized
153 that positive associations between defense systems and MGE abundance could be explained, at least
154 in part, by recent co-transfer events (henceforth, we term this the “linkage hypothesis”). Under this
155 hypothesis, positive associations would simply result from defense systems travelling together with
156 MGE.

157 To test the linkage hypothesis, we first verified that defense systems tend to be co-transferred with
158 MGE. We used ancestral reconstruction methods to identify the branches in which defense systems
159 and MGE were gained and lost along strain trees (Table S5). We found that >95% of defense
160 acquisition events occurred in branches in which MGE were also gained, and >90% of defense losses

161 occurred in branches in which MGE were also lost (Fig. 4A). The random expectation given the rates
162 of MGE gain and loss would be 71% and 63%, respectively (deviations with respect to these
163 expectations are statistically significant with $p < 10^{-8}$, binomial test). We also quantified the effect of
164 MGE gain and loss on the per-branch probability of gaining or losing defense systems. The probability
165 of acquiring a defense system is around 50-fold higher in branches in which MGE are gained than in
166 other branches (Fig. 4B-C; Fisher's exact test $p < 10^{-8}$ for all defense systems). Similarly, the
167 probability of losing a defense system is 10- to 50-fold higher if MGE are also lost in the same branch
168 (Fig. 4B-C; Fisher's exact test $p < 10^{-20}$ for all defense systems). These trends are observed even in
169 very short branches, spanning an evolutionary time of 10^{-6} substitutions per site in core genes (Fig.
170 S5). To rule out the possibility that these associations were due to the presence of incomplete
171 genomes, we separately considered terminal and internal tree branches. Because the dataset only
172 includes complete and nearly complete genomes, the absence of a small number of missing genes, if
173 relevant, would only affect the inference of gene gain and loss in terminal branches (those leading
174 from the immediate ancestor to each incomplete genome). Despite some quantitative differences,
175 strong co-acquisition (and co-loss) of MGE and defense occurs in both terminal and internal branches
176 (Fig. S5), confirming that these associations are genuine.

177 Although co-gain of MGE and defense systems along the same tree branch does not necessarily imply
178 a single event of joint gain, the strength of the association, the extremely short timespans, and the
179 fact that similar trends are observed for gene losses strongly suggest that concurrent gain and loss of
180 MGE and defense systems involve genetic linkage. (Alternative explanations in terms of strong
181 selective pressure to quickly acquire defense mechanisms upon exposure to MGE and lose them
182 once the MGE disappear might produce similar trends at the population level but not at the single-
183 genome level, whereas episodic increases in the overall rate of HGT leading to separate but
184 correlated acquisition of MGE and defense systems would not explain correlated losses.)

185 A major consequence of the co-transfer of MGE and defense systems is that the negative effects of
186 defense systems on HGT should be easier to detect at longer timescales or under evolutionary
187 conditions that weakened genetic linkage. To test that prediction, we studied how relative
188 differences in the rates of gene gain between DEF⁺ and DEF⁻ clades depend on the depth of their last
189 common ancestor in the strain tree (note that the depth of the last common ancestor serves as an
190 upper limit for the time that the defense system has been retained in a lineage). As expected,
191 positive outliers (with much higher gene gain rate in the DEF⁺ clade than in the DEF⁻ clade) almost
192 invariably correspond to very recent lineages (less than 10⁻⁵ substitutions per bp in nearly universal
193 core genes), indicating a very recent acquisition of the defense system (Fig. 5A and Table S6). In
194 contrast, negative outliers (with much lower gene gain rate in the DEF⁺ clade than in the DEF⁻ clade)
195 generally correspond to deeper lineages (at least 0.001 substitutions per bp). We quantified these
196 trends by calculating the skewness of the distribution of DEF⁺/DEF⁻ log-transformed gain ratios in
197 very recent and older lineages (Fig. 5B and Table S7). All the distributions show significant positive
198 skewness in recent lineages and significant negative skewness in older lineages (all $p < 0.001$ except
199 recent RM lineages, with $p = 0.023$; d'Agostino test for skewness). This quantitative analysis confirms
200 that differences in timescale affect all defense systems, with positive associations between defense
201 systems and gene gain being dominant in the short term and negative associations becoming more
202 frequent in the long term.

203 A second prediction of the linkage hypothesis is that the net effect of defense systems on HGT is not
204 only species specific, but also lineage specific. That is, defense systems may display positive, zero, or
205 negative association with HGT in different lineages depending on when the defense system was
206 acquired and how tight is the linkage to MGE. A more detailed study of gene gain rates in individual
207 species confirms that the effect of CRISPR-Cas is, indeed, lineage-specific, with the same species
208 encompassing recent CRISPR-Cas⁺ lineages that display increased gene gain rates and older CRISPR-
209 Cas⁺ lineages with reduced gene gain rates (Fig. 5C). This observation, combined with the very recent
210 acquisition of CRISPR-Cas in most lineages, could explain why many species show nonsignificant or

211 positive associations between CRISPR-Cas and MGE abundances. Indeed, of the six representative
212 species shown in Fig. 5C, only *Streptococcus anginosus* and *Acinetobacter radioresistens* show a net
213 negative association between CRISPR-Cas and MGE abundance in the PGLMM analysis.

214 The findings described so far indicate that, of the seven defense systems considered in this study,
215 CRISPR-Cas is the one that most often induces a net reduction in HGT rates. According to the linkage
216 hypothesis, this could be due to a comparatively weaker physical association between CRISPR-Cas
217 and MGE. To evaluate that possibility, we calculated the percentage of genes from CRISPR-Cas
218 located inside MGE and compared that to other defense systems (Fig. 5D). Consistent with the
219 linkage hypothesis, CRISPR-Cas is the defense system that is least often encoded by MGE (10.0% vs
220 20.1% for all other systems; $p < 10^{-10}$, Fisher's exact test), followed by RM (16.7%) and DMS (17.0%).

221 Anti-CRISPR proteins modulate associations between CRISPR-Cas and HGT

222 Anti-CRISPR proteins (Acr) have the potential to suppress the possible negative effect of CRISPR-Cas
223 on MGE-driven gene transfer (Mahendra et al. 2020). This, in turn, could contribute to explaining the
224 sign of the association between CRISPR-Cas and MGE in different lineages. To assess that possibility,
225 we searched all the genomes in the dataset for known Acr, finding 16,058 proteins. Then, we
226 compared the prevalence of Acr in species in which the presence of CRISPR-Cas is positively or
227 negatively correlated with the MGE content, separately considering genomes that do and do not
228 harbor CRISPR-Cas. Because Acr are generally encoded by MGE (Pinilla-Redondo et al. 2020) and the
229 prevalence of MGE systematically varies across groups acting as a confounding factor, we restricted
230 our comparisons to genomes that contain at least one prophage. Our results indicate that genomes
231 with and without CRISPR-Cas differ in the prevalence of Acr (Fig. S6). More importantly, these
232 differences have opposite directions in species in which CRISPR-Cas is positively and negatively
233 associated with MGE abundance. In the former, Acr are slightly more prevalent in genomes that
234 contain CRISPR-Cas (46% vs 43%, $p = 0.026$, Fisher's exact test). In the latter, Acr are less prevalent in
235 genomes with CRISPR-Cas (35% vs 41%, $p < 10^{-4}$, Fisher's exact test). In both groups, the prevalence

236 of Acr in genomes without CRISPR-Cas is similar. These results suggest that negative associations
237 between CRISPR-Cas and HGT could be dependent on (or at least facilitated by) low prevalence of Acr
238 in the genome.

239

240 **Discussion**

241 Defense systems could have a significant impact on microbial evolution by effectively blocking the
242 transfer of MGE, reducing gene flow and limiting the spread of accessory genes. However, because
243 defense systems are often carried by MGE and these are the main vehicles of HGT, a net positive
244 association between defense systems and gene exchange cannot be ruled out *a priori*. We assessed
245 the relative weight of these two opposite scenarios by quantifying the association between defense
246 systems, MGE abundance, and gene acquisition rates in a phylogeny-aware comparative study of 197
247 prokaryotic species.

248 Our results shed light on previous, apparently contradictory findings concerning the effect of CRISPR-
249 Cas on genome evolution and diversification (Gophna et al. 2015; O'Meara and Nunney 2019;
250 Shehreen et al. 2019; Wheatley and MacLean 2021; Pursey et al. 2022). A pioneering study
251 conducted in 2015 found no evidence to support an overall association between CRISPR-Cas activity
252 and reduced gene acquisition via HGT at evolutionary timescales (Gophna et al. 2015). Such lack of
253 association was explained by several factors, including the high mobility of CRISPR-Cas systems, that
254 limits their long-term impact on host genomes, and the possibility that HGT is mediated by MGE that
255 escape (or are not targeted by) CRISPR-Cas immunity.

256 Our analyses support the general conclusion that CRISPR-Cas and other defense systems have little
257 overall impact on HGT in most bacterial species. Specifically, differences in the rates of gene
258 acquisition in lineages that do and do not harbor defense systems are centered around zero. That
259 said, we identified significant opposite trends at very short and intermediate evolutionary timescales:
260 positive associations between defense systems and HGT are more frequent at very short timescales,

261 whereas negative associations become dominant at longer timescales. These opposite trends suggest
262 that the possible negative effects of defense systems on gene exchange may be obscured by recent
263 co-transfer events involving MGE. Indeed, we found strong evidence of linkage between MGE and
264 defense systems, which tend to be not only co-acquired but also co-lost from bacterial genomes at
265 fast rates. Such joint dynamics of MGE and defense systems manifest as a pattern of positive
266 associations at short timescales. In turn, the negative effects of anti-MGE defense on HGT only
267 become detectable if defense systems are maintained for long enough periods of time in the
268 chromosome or in stable plasmids.

269 Besides this general picture, we identified some species in which the presence of defense systems
270 (especially CRISPR-Cas) significantly correlates with smaller genome sizes and MGE abundances.
271 Some of those associations had been previously described in *P. aeruginosa* and *Klebsiella*
272 *pneumoniae* (Wheatley and MacLean 2021; Botelho et al. 2023). And yet, these species represent
273 special cases rather than the rule, even in the context of host-associated bacteria. In fact, our results
274 underline that the association between defense systems and HGT is strongly system- and lineage-
275 dependent. This conclusion confirms and extends previous findings that showed that the impact of
276 CRISPR-Cas on the spread of antibiotic resistance is highly variable across species and its sign cannot
277 be easily explained by simple ecological, environmental, or genomic variables (Shehreen et al. 2019).

278 The statistical modeling approach used in this study implies several assumptions and limitations that
279 are worth discussing. First, phylogenetic corrections in PGLMM assume a linear association between
280 the covariate (measured as nucleotide divergence in core genes) and the response variable (gene
281 abundances). Therefore, the model can accommodate cross-strain variation in evolutionary rates as
282 long as such variation also affects, in a correlated manner, the rates of gene gain and loss. Otherwise,
283 heterogeneity in evolutionary rates would reduce the explanatory power of the phylogenetic term
284 and the sensitivity of statistical tests. Second, variability in sample size, gene abundance, and
285 prevalence of both defense systems and MGE results in heterogeneous statistical power across

286 species and systems. We introduced an effect size criterion to minimize possible biases, but one
287 should still be cautious when interpreting some of the results. In particular, absolute numbers of
288 positive and negative associations may be subject to sensitivity biases, whereas trends regarding the
289 balance of positive and negative associations are likely more robust. Finally, we fitted all RM and
290 DMS with one single model, even if the restriction sites and epigenetic modifications of these
291 systems are very diverse. This could lead to an underestimation of effect sizes and statistical
292 significance due to the clumping of dissimilar trends.

293 Among the 7 defense systems included in this study, CRISPR-Cas stands out for being more often
294 associated with reduced genome sizes and lower MGE (especially prophage) abundances. In contrast,
295 DMS, Abi, DRT, and CBASS are more often associated with higher numbers of MGE and accessory
296 genes. This finding is fully consistent with a recent study that compared 73 defense systems in 12
297 bacterial species (Kogay et al. 2024). We propose that what makes CRISPR-Cas systems different is
298 their weaker (though still substantial) linkage with MGE. Compared to fully functional CRISPR-Cas
299 systems, other defense systems like Abi, Gabija, CBASS, DMS, and RM are more frequently located
300 within or next to MGE (Makarova et al. 2011; Benler et al. 2021; Rousset et al. 2022; Botelho 2023)
301 and may have alternative functions related to MGE propagation. For example, Abi systems have been
302 identified in PICIs as accessory genes that facilitate their parasitic lifecycle (Ibarra-Chavez et al. 2021).
303 The narrow specificity of some defense systems may be another reason why those systems do not
304 significantly interfere with HGT. For instance, the GmrSD type IV RM system selectively targets
305 phages with glucosylated hydroxymethylcytosine (Bair and Black 2007) and the *Thoeris* defense
306 system only appears to be effective against myoviruses (Doron et al. 2018). This caveat extends to
307 any defense system based on epigenetic modifications, such as RM and DMS, whose overall effect on
308 HGT depend on the repertoire of epigenetic markers in the host and MGE populations (Oliveira et al.
309 2016). Although selective targeting is often viewed as an outcome of phage-host coevolution, it is
310 tempting to speculate that it could have been evolutionarily favored by the need to fight harmful

311 genetic parasites while maintaining sufficiently high rates of HGT to prevent population-level gene
312 loss (Iranzo et al. 2016).

313 Defense systems sometimes exhibit synergistic interactions (Dupuis et al. 2013; Wu et al. 2024),
314 which could contribute to the heterogeneity of effects reported in this study. The vast number of
315 potential interactions and the limited availability of high-quality genomes made it unfeasible to
316 systematically account for the effect of interactions with the methods developed in this work.
317 Determining if and to what extent synergy among defense systems affects HGT remains a subject for
318 future investigation, possibly focused on a small set of experimentally validated interactions in highly
319 sequenced species. Another open question concerns which levels of detail, both taxonomic and
320 functional, best capture the effect of defense systems on HGT. From a taxonomic perspective,
321 working at or below the species level is a natural choice because species represent genetically
322 cohesive units (Bobay and Ochman 2017; Konstantinidis 2023; Conrad et al. 2024) and, as a result,
323 uncontrolled confounding factors are less likely to affect within-species than cross-species
324 comparisons. In contrast, more complex multi-level approaches would be required to detect trends
325 at higher taxonomic ranks. From a functional perspective, we grouped defense systems based on
326 their mechanism of action, under the assumption that functionally similar systems produce similar
327 effects on HGT. Though reasonable, this grouping criterion may not be optimal in systems in which
328 subtypes markedly differ in their eco-evolutionary dynamics and linkage with MGE. Moreover, fine-
329 grain dissection of highly abundant systems, such as RM and DMS, could help improve the sensitivity
330 of statistical tests by producing more balanced sets of strains with and without the subtypes of
331 interest.

332 All in all, we showed that some defense systems, especially CRISPR-Cas, can significantly reduce HGT,
333 although the effect is often masked by the fact that these systems travel together with MGE. Beyond
334 possible functional connections, the linkage between defense systems and MGE is an inevitable
335 consequence of the arms race between parasites and hosts. Because defense systems are costly and

336 their efficacy drops as parasites evolve, they are subject to rapid turnover and depend on HGT for
337 long-term persistence in microbial populations (van Houte et al. 2016; Iranzo et al. 2017; Koonin et
338 al. 2017; Puigbo et al. 2017). As a result, it is extremely challenging to disentangle the impact of
339 defense systems on gene flow from the causes that lead to their presence or absence, especially at
340 short evolutionary timescales. As more and more genomic data become available, we expect that
341 future research will overcome this challenge by quantifying the linkage between defense systems
342 and MGE, developing more realistic null models, and testing the role of defense systems on microbial
343 adaptation at different timescales.

344

345 **Methods**

346 Genome collection and identification of defense systems

347 We parsed the Genome Taxonomy Database (GTDB, <https://gtdb.ecogenomic.org>) release 202 (Parks
348 et al. 2020) to identify all high-quality genomes (according to the MIMAG criteria (Bowers et al.
349 2017)) with completeness >99%, contamination <1%, and contig count <500. The 82,595 genomes
350 that passed these filters were downloaded from the NCBI FTP site (<https://ftp.ncbi.nlm.nih.gov>).
351 CRISPR-Cas systems were identified with CRISPRCasTyper v1.2.4 (Russel et al. 2020) using default
352 parameters. We classified a genome as CRISPR-Cas⁺ if it contains at least one high-confidence Cas
353 operon and CRISPR-Cas⁻ otherwise. Only the species with >10 genomes and at least 5 CRISPR-Cas⁺
354 genomes were further considered. To reduce the computational cost, we only considered a
355 maximum of 500 genomes per species. Species with >500 genomes were randomly subsampled to
356 keep at most 350 CRISPR-Cas⁺ and 150 CRISPR-Cas⁻ genomes. Other defense systems were identified
357 with Padloc v1.1.0 (db v1.4.0) (Payne et al. 2022). Our analysis focused on the most prevalent
358 defense systems: restriction-modification (RM), DMS, Abi, CRISPR-Cas, Gabija, DRT, and CBASS.

359 After applying these criteria, 19,323 genomes belonging to 196 bacterial and 1 archaeal species
360 (*sensu* GTDB) were included in the analysis (Tables S1 and S2). Of those, 2,964 correspond to
361 complete genomes and the rest to high-quality, nearly complete ones.

362 Gene prediction and annotation.

363 Open reading frames (ORF) were predicted with Prodigal v2.6.3, using codon table 11 (prokaryotic
364 genetic code) and “single” mode, as recommended for finished and draft quality genomes (Hyatt et
365 al. 2010). Orthologous ORF were then separately clustered for each species with Roary v3.13.0 (Page
366 et al. 2015) setting an 80% identity threshold for initial clustering followed by synteny-based
367 refinement (options ‘-t 11 -i 80’). The resulting gene clusters were functionally annotated by selecting
368 a representative sequence, arbitrarily chosen among those with length between 0.95 and 1.05 times
369 the average length of all sequences in the cluster. Representative sequences were functionally
370 annotated by mapping them to in-house profiles of the Clusters of Orthologous Genes (COG)
371 database (2020 release) (Galperin et al. 2021) with HMMER v3.1b2 (e-value < 0.001)
372 (<http://hmmer.org>). The 26 major prokaryotic functional categories defined in the COG database
373 were assigned to the annotated genes. Some functional categories (A, RNA processing and
374 modification; B, chromatin structure and dynamics; W, extracellular structures; T, signal
375 transduction; and Z, cytoskeleton) were excluded since they rarely or never occur in prokaryotic
376 genomes. The case-insensitive keywords “phage”, “plasmid” and “transpos*” in the COG gene
377 annotations were used to identify genes associated with prophages, plasmids, and transposons,
378 respectively, and the resulting gene lists were manually curated to minimize false assignments. We
379 based our statistical analyses on marker gene counts rather than full MGE counts because the latter
380 are more susceptible to technical artifacts (e.g., different heuristics to deal with nested MGE will
381 affect the number of MGE, but not the number of MGE marker genes). Gene counts per genome and
382 functional category are listed in Table S8. Anti-defense proteins, including anti-CRISPR, were
383 identified by running HMMER v3.1b2 (e-value < 10^{-10}) against the dbAPIS database (Yan et al. 2024).

384 Identification of genomic regions containing mobile genetic elements

385 Prophages and plasmids were detected with geNomad v1.8.1 using default options (Camargo et al.
386 2024). Short transposons were identified based on the presence of isolated or paired genes
387 annotated as transposases. In the latter case, we allowed for up to one additional gene between two
388 transposon-related genes to account for the genetic architecture of some insertion sequences
389 (Gomez et al. 2014). ICEfinder (Liu et al. 2019) was used to identify ICE and IME. The MGE identified
390 through these approaches were masked from complete genomes to produce the gene counts in
391 Table S9 and the results shown in Fig S2.

392 Species trees

393 Phylogenetic trees were separately built for each species based on the set of 120 prokaryotic marker
394 genes (122 in the case of Archaea) proposed by the GTDB r202. For each species, only those marker
395 genes with prevalence >80% were used for phylogenetic reconstruction. We aligned the amino acid
396 sequences of each marker gene with mafft-linsi (L-INS-I algorithm, default options, MAFFT v7.475)
397 (Katoh and Standley 2013) and back-translated the amino acid alignments to nucleotide alignments
398 with pal2nal.pl v14 (Suyama et al. 2006) using codon table 11. After concatenating all nucleotide
399 alignments, we built preliminary trees with FastTree v2.1.10 (options '-gtr -nt -gamma -nosupport -
400 mlacc 2 -slownni') (Price et al. 2010). The tree topologies produced by FastTree were subsequently
401 provided to RAxML v8.2.12 (Stamatakis 2014) for branch length optimization (raxmlHPC with options
402 '-f e -m GTRGAMMA'). The final trees are included in Supplementary File S1.

403 To visualize trends across species (Fig. 2, S3, and S4), we used the online tool iTOL (Letunic and Bork
404 2021) and the multispecies tree from GTDB r202.

405 Phylogenetic generalized linear mixed models (PGLMM)

406 For each genome in the dataset, we collected the following response variables: the total number of
407 genes, the number of genes belonging to each functional category, and the number of genes

408 associated with prophages, plasmids, and transposons (Table S8). Genes that belong to the 7 defense
409 systems of interest were excluded when computing these values. Then, for each response variable,
410 we fitted a PGLMM with Poisson distribution and canonical link function, using the presence or
411 absence of each defense system as predictors and the species trees as guides to generate the
412 covariance matrix.

413 For each species, the PGLMM assumes that the response variable, Y_i , follows a Poisson distribution
414 with mean μ_i , that is, $Y_i \sim \text{Poisson}(\mu_i)$. The expected gene abundance, μ_i , is modeled as
415 $\log \mu_i = \beta_0 + \sum \beta_j X_{ij} + \epsilon_i$, where β_0 is the intercept, β_j is the coefficient associated with the defense
416 system j , and $X_{ij} \in \{0,1\}$ denotes the absence or presence of defense system j in genome i . The
417 random effects ϵ_i follow a multivariate normal distribution, $\epsilon \sim \text{Gaussian}(0, \sigma_{phy}^2 \mathbf{C})$, where σ_{phy}^2 is the
418 strength of the phylogenetic signal and \mathbf{C} is a covariance matrix derived from the phylogenetic tree
419 under the assumption of Brownian motion evolution.

420 To extend the PGLMM to multiple species, we considered that the effect of defense systems on gene
421 content may not be the same across species (this was a reasonable assumption *a priori* and later
422 confirmed by the analysis). To account for that, the multispecies model must include an interaction
423 term “defense \times species”, whose coefficients are relevant per se, and, accordingly, modeled as fixed
424 effects. Moreover, because HGT rates are highly variable among species (Puigbò et al. 2014; Iranzo et
425 al. 2019), it would be inadequate to extrapolate phylogenetic corrections beyond single species or
426 assume that the strength of the phylogenetic signal is the same for all species. These considerations
427 are captured by a phylogenetic covariance matrix with block-diagonal structure (one block per
428 species), and species-wise values of the phylogenetic coefficient (one for each block of the
429 phylogenetic covariance matrix). In practice, fitting a multispecies model with these specifications is
430 formally equivalent to fitting independent models for each species. The latter approach, that we
431 adopted, has the advantage of being more suitable for parallelization and requiring fewer
432 computational resources. Thus, for each species and response variable, we fitted a PGLMM with the

433 function `pglmm_compare(response_variable ~ Abi + CBASS + CRISPR + DMS + DRT + Gabija + RM,`
434 `family = "poisson", data = SpData, phy = SpTree)` from the R package `phyr v1.1.0` (Li et al. 2020). In
435 the formula, `Abi`, `CBASS`, `CRISPR`, and so on, are binary variables representing the presence (1) or
436 absence (0) of each defense system. Table S4 presents the coefficients, p-values, and goodness of fit
437 of the model.

438 The model described above includes nine coefficients (intercept, seven defense systems, and the
439 phylogenetic signal), which could lead to overfitting in species in which the number of available
440 genomes is limited. As an alternative, we also fitted seven separate PGLMM, one for each defense
441 system, involving a single predictor and the phylogenetic random effect (Table S3). This approach
442 does not account for correlations among defense systems, but it has the advantage of not being
443 affected by overfitting. The figures in the manuscript are based on this second set of models,
444 although, in practice, both approaches produce very similar quantitative results.

445 For each class of PGLMM, we also fitted non-phylogenetic models in which the covariance matrix **C** of
446 the random effect was replaced by the identity matrix. The PGLMM were compared to their non-
447 phylogenetic counterparts using the conditional Akaike Information Criterion (cAIC) as previously
448 described (Greven and Kneib 2010; Säfken et al. 2021). Based on the cAIC, PGLMM performed better
449 than non-phylogenetic models in 65% of the species-response-defense triplets and were close to
450 non-phylogenetic models ($\Delta cAIC < 2$) in another 25% of the triplets (in most of those cases, the
451 strength of the phylogenetic signal was close to zero, which made both phylogenetic and non-
452 phylogenetic models equivalent).

453 To account for the possibility that other (less abundant) defense systems could explain part of the
454 variability in the results, we explored a more complex set of models that included the total number
455 of other defense systems as an additional predictor. These models generally performed worse than
456 their simpler variants ($\Delta cAIC > 0$ in 83% of the species) and were not considered for further analysis.

457 Large differences in sample size among species and defense systems translate into unequal precision
458 in the estimation of the model coefficients. To deal with that limitation, we used two different
459 criteria to identify species in which the presence of a defense system is positively or negatively
460 associated with gene numbers. For one option, we adopted a classical criterion of statistical
461 significance ($p < 0.05$) for the predictor variable in the PGLMM. For the other option, we applied an
462 alternative criterion based on effect size, aimed at comparing species with different sample sizes in
463 which p-values are not commensurable. Specifically, for each variable and defense system, we jointly
464 considered the PGLMM of all the species and determined the smallest effect size (in absolute value)
465 that reached statistical significance ($p < 0.05$) in any species. Then we used the smallest significant
466 effect size (SSES) as a threshold to classify associations as positive, negative, or null.

467 Inference of gene gain and loss

468 Gene gains and losses at each branch of each species tree were estimated with Gloome (Cohen and
469 Pupko 2010), using as inputs the gene presence/absence matrices previously generated by Roary and
470 the species trees. The parameter configuration file was set to optimize the likelihood of the observed
471 phyletic profiles under a genome evolution model with 4 categories of gamma-distributed gain and
472 loss rates and stationary frequencies at the root.

473 Comparison of gene gain and loss rates between DEF⁺ and DEF⁻ clades

474 For each species tree, we defined DEF⁺ clades as the narrowest possible clades such that at least 80%
475 of the leaves contain the defense system of interest. Candidate DEF⁻ clades were defined in an
476 analogous way, referring to leaves without the defense system. Next, we identified pairs of DEF⁺ and
477 DEF⁻ clades that constitute sister groups. Sister DEF^{+/-} pairs were excluded if both clades contained a
478 single genome. For each clade in a valid pair, we computed the overall gene gain and loss rates as the
479 expected number of gene gains (or losses) in that clade divided by the total branch length. Gain and
480 loss rates for different functional categories and MGE we calculated in an analogous way but
481 restricting the sum of gene gains and losses to the genes of interest. When calculating overall and

482 category-wise gain and loss rates, we did not take into account the contribution of species-wise
483 singletons (genes without homologs in other genomes of the same species), as they may represent
484 false gene predictions or genes that are replaced at unusually high rates (Wolf et al. 2016). To
485 account for the non-negative nature of gain and loss rates and their heavy-tailed distributions,
486 comparisons between DEF^+ and DEF^- sister branches were done based on log-transformed rate
487 estimates. To calculate the skewness of the distributions and their statistical significance, we use the
488 method proposed by D'Agostino (D'Agostino et al. 1990) as implemented by the
489 `scipy.stats.skewtest()` function in Python (Virtanen et al. 2020).

490

491 **Competing interests**

492 The authors declare no competing financial interests.

493 **Acknowledgements**

494 Y.L. is supported by China Scholarship Council (No.202008440425). J.B. is supported by the Maria
495 Zambrano grant of the Spanish Ministry of Universities (Grant No. UP2021-035), and the Severo
496 Ochoa Program for Centres of Excellence in R&D of the Agencia Estatal de Investigación of Spain
497 (Grant No. CEX2020-000999-S (2022–2025) to the CBGP). J.I. is supported by the Ramón y Cajal
498 Programme of the Spanish Ministry of Science (Grant No. RYC-2017–22524); the Agencia Estatal de
499 Investigación of Spain (Grant Nos. PID2019-106618GA-I00 and CNS2023-145430), the Severo Ochoa
500 Programme for Centres of Excellence in R&D of the Agencia Estatal de Investigación of Spain (Grant
501 No. SEV-2016–0672 (2017–2021) to the CBGP); and the Comunidad de Madrid (through the call
502 Research Grants for Young Investigators from Universidad Politécnica de Madrid, Grant No.
503 M190020074JIIS).

504 We thank Jorge Calle-Espinosa for helpful discussions and Lindsay Dudbridge for critical reading of
505 the manuscript.

506 Author contributions: Y.L. conducted bioinformatic and statistical analyses; J.B. co-supervised the
 507 study; J.I. designed and supervised the study. All authors contributed to the interpretation of the
 508 results and writing of the manuscript.

509

510 References

- 511 Arnold BJ, Huang IT, Hanage WP. 2022. Horizontal gene transfer and adaptive evolution in bacteria. *Nat Rev*
 512 *Microbiol* **20**: 206-218.
- 513 Bair CL, Black LW. 2007. A type IV modification dependent restriction nuclease that targets glucosylated
 514 hydroxymethyl cytosine modified DNAs. *J Mol Biol* **366**: 768-778.
- 515 Benler S, Faure G, Altae-Tran H, Shmakov S, Zheng F, Koonin E. 2021. Cargo Genes of Tn7-Like Transposons
 516 Comprise an Enormous Diversity of Defense Systems, Mobile Genetic Elements, and Antibiotic
 517 Resistance Genes. *mBio* **12**: e0293821.
- 518 Bobay LM, Ochman H. 2017. Biological species are universal across life's domains. *Genome biology and*
 519 *evolution* **9**: 491-501.
- 520 Botelho J. 2023. Defense systems are pervasive across chromosomally integrated mobile genetic elements and
 521 are inversely correlated to virulence and antimicrobial resistance. *Nucleic acids research* **51**: 4385-
 522 4397.
- 523 Botelho J, Cazares A, Schulenburg H. 2023. The ESKAPE mobilome contributes to the spread of antimicrobial
 524 resistance and CRISPR-mediated conflict between mobile genetic elements. *Nucleic acids research* **51**:
 525 236-252.
- 526 Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, Schulz F, Jarett J, Rivers AR,
 527 Eloe-Fadrosh EA et al. 2017. Minimum information about a single amplified genome (MISAG) and a
 528 metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nature biotechnology* **35**: 725-
 529 731.
- 530 Camargo AP, Roux S, Schulz F, Babinski M, Xu Y, Hu B, Chain PSG, Nayfach S, Kyrpides NC. 2024. Identification of
 531 mobile genetic elements with geNomad. *Nature biotechnology* **42**: 1303-1312.
- 532 Cohen O, Pupko T. 2010. Inference and characterization of horizontally transferred gene families using
 533 stochastic mapping. *Molecular biology and evolution* **27**: 703-713.
- 534 Conrad RE, Brink CE, Viver T, Rodriguez-R LM, Aldeguer-Riquelme B, Hatt JK, Venter SN, Amann R, Rossello-
 535 Mora R, Konstantinidis KT. 2024. Microbial species exist and are maintained by ecological
 536 cohesiveness coupled to high homologous recombination. *bioRxiv* doi:10.1101/2024.05.25.595874:
 537 2024.2005.2025.595874.
- 538 Costa AR, van den Berg DF, Esser JQ, Muralidharan A, van den Bossche H, Bonilla BE, van der Steen BA,
 539 Haagsma AC, Fluit AC, Nobrega FL et al. 2024. Accumulation of defense systems in phage-resistant
 540 strains of *Pseudomonas aeruginosa*. *Sci Adv* **10**: eadj0341.
- 541 D'Agostino RB, Belanger A, D'Agostino Jr RB. 1990. A Suggestion for Using Powerful and Informative Tests of
 542 Normality. *The American Statistician* **44**: 316-321.
- 543 Doron S, Melamed S, Ofir G, Leavitt A, Lopatina A, Keren M, Amitai G, Sorek R. 2018. Systematic discovery of
 544 antiphage defense systems in the microbial pangenome. *Science* **359**.
- 545 Dupuis ME, Villion M, Magadan AH, Moineau S. 2013. CRISPR-Cas and restriction-modification systems are
 546 compatible and increase phage resistance. *Nature communications* **4**: 2087.
- 547 Galperin MY, Wolf YI, Makarova KS, Vera Alvarez R, Landsman D, Koonin EV. 2021. COG database update: focus
 548 on microbial diversity, model organisms, and widespread pathogens. *Nucleic acids research* **49**: D274-
 549 D281.
- 550 Georjon H, Bernheim A. 2023. The highly diverse antiphage defence systems of bacteria. *Nat Rev Microbiol* **21**:
 551 686-700.
- 552 Gomez MJ, Diaz-Maldonado H, Gonzalez-Tortuero E, Lopez de Saro FJ. 2014. Chromosomal replication
 553 dynamics and interaction with the beta sliding clamp determine orientation of bacterial transposable
 554 elements. *Genome biology and evolution* **6**: 727-740.

- 555 Gophna U, Kristensen DM, Wolf YI, Popa O, Drevet C, Koonin EV. 2015. No evidence of inhibition of horizontal
556 gene transfer by CRISPR-Cas on evolutionary timescales. *The ISME journal* **9**: 2021-2027.
- 557 Greven S, Kneib T. 2010. On the behaviour of marginal and conditional AIC in linear mixed models. *Biometrika*
558 **97**: 773-789.
- 559 Hall RJ, Whelan FJ, McInerney JO, Ou Y, Domingo-Sananes MR. 2020. Horizontal Gene Transfer as a Source of
560 Conflict and Cooperation in Prokaryotes. *Front Microbiol* **11**: 1569.
- 561 Haudiquet M, de Sousa JM, Touchon M, Rocha EPC. 2022. Selfish, promiscuous and sometimes useful: how
562 mobile genetic elements drive horizontal gene transfer in microbial populations. *Philos Trans R Soc*
563 *Lond B Biol Sci* **377**: 20210234.
- 564 Hyatt D, Chen GL, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal: prokaryotic gene recognition
565 and translation initiation site identification. *BMC Bioinformatics* **11**: 119.
- 566 Ibarra-Chavez R, Hansen MF, Pinilla-Redondo R, Seed KD, Trivedi U. 2021. Phage satellites and their emerging
567 applications in biotechnology. *FEMS Microbiol Rev* **45**.
- 568 Irazo J, Cuesta JA, Manrubia S, Katsnelson MI, Koonin EV. 2017. Disentangling the effects of selection and loss
569 bias on gene dynamics. *Proceedings of the National Academy of Sciences of the United States of*
570 *America* **114**: E5616-E5624.
- 571 Irazo J, Puigbo P, Lobkovsky AE, Wolf YI, Koonin EV. 2016. Inevitability of Genetic Parasites. *Genome biology*
572 *and evolution* **8**: 2856-2869.
- 573 Irazo J, Wolf YI, Koonin EV, Sela I. 2019. Gene gain and loss push prokaryotes beyond the homologous
574 recombination barrier and accelerate genome sequence divergence. *Nature communications* **10**: 5376.
- 575 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in
576 performance and usability. *Molecular biology and evolution* **30**: 772-780.
- 577 Kogay R, Wolf YI, Koonin EV. 2024. Defence systems and horizontal gene transfer in bacteria. *Environmental*
578 *Microbiology* **26**: e16630.
- 579 Konstantinidis KT. 2023. Sequence-discrete species for prokaryotes and other microbes: A historical perspective
580 and pending issues. *mLife* **2**: 341-349.
- 581 Koonin EV, Makarova KS, Wolf YI. 2017. Evolutionary Genomics of Defense Systems in Archaea and Bacteria.
582 *Annual review of microbiology* **71**: 233-261.
- 583 Koonin EV, Makarova KS, Wolf YI, Krupovic M. 2020. Evolutionary entanglement of mobile genetic elements
584 and host defence systems: guns for hire. *Nature reviews Genetics* **21**: 119-131.
- 585 Lee IPA, Eldakar OT, Gogarten JP, Andam CP. 2022. Bacterial cooperation through horizontal gene transfer.
586 *Trends Ecol Evol* **37**: 223-232.
- 587 Letunic I, Bork P. 2021. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display and
588 annotation. *Nucleic acids research* **49**: W293-W296.
- 589 Li D, Dinnage R, Nell LA, Helmus MR, Ives AR. 2020. phyr: An R package for phylogenetic species-distribution
590 modelling in ecological communities. *Methods in Ecology and Evolution* **11**: 1455-1463.
- 591 Liu M, Li X, Xie Y, Bi D, Sun J, Li J, Tai C, Deng Z, Ou HY. 2019. ICEberg 2.0: an updated database of bacterial
592 integrative and conjugative elements. *Nucleic acids research* **47**: D660-D665.
- 593 Mahendra C, Christie KA, Osuna BA, Pinilla-Redondo R, Kleinstiver BP, Bondy-Denomy J. 2020. Broad-spectrum
594 anti-CRISPR proteins facilitate horizontal gene transfer. *Nature microbiology* **5**: 620-629.
- 595 Makarova KS, Wolf YI, Snir S, Koonin EV. 2011. Defense islands in bacterial and archaeal genomes and
596 prediction of novel defense systems. *Journal of bacteriology* **193**: 6039-6056.
- 597 Marraffini LA, Sontheimer EJ. 2008. CRISPR interference limits horizontal gene transfer in staphylococci by
598 targeting DNA. *Science* **322**: 1843-1845.
- 599 Mayo-Munoz D, Pinilla-Redondo R, Birkholz N, Fineran PC. 2023. A host of armor: Prokaryotic immune
600 strategies against mobile genetic elements. *Cell Rep* **42**: 112672.
- 601 Meaden S, Biswas A, Arkhipova K, Morales SE, Dutilh BE, Westra ER, Fineran PC. 2022. High viral abundance
602 and low diversity are associated with increased CRISPR-Cas prevalence across microbial ecosystems.
603 *Curr Biol* **32**: 220-227 e225.
- 604 Morris JA, Gardner MJ. 1988. Calculating confidence intervals for relative risks (odds ratios) and standardised
605 ratios and rates. *Br Med J (Clin Res Ed)* **296**: 1313-1316.
- 606 O'Hara BJ, Barth ZK, McKitterick AC, Seed KD. 2017. A highly specific phage defense system is a conserved
607 feature of the *Vibrio cholerae* mobilome. *PLoS genetics* **13**: e1006838.
- 608 O'Meara D, Nunnery L. 2019. A phylogenetic test of the role of CRISPR-Cas in limiting plasmid acquisition and
609 prophage integration in bacteria. *Plasmid* **104**: 102418.

- 610 Oliveira PH, Touchon M, Rocha EP. 2016. Regulation of genetic flux between bacteria by restriction-
611 modification systems. *Proceedings of the National Academy of Sciences of the United States of*
612 *America* **113**: 5658-5663.
- 613 Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D, Keane JA, Parkhill J. 2015.
614 Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**: 3691-3693.
- 615 Parks DH, Chuvochina M, Chaumeil PA, Rinke C, Mussig AJ, Hugenholtz P. 2020. A complete domain-to-species
616 taxonomy for Bacteria and Archaea. *Nature biotechnology* **38**: 1079-1086.
- 617 Payne LJ, Meaden S, Mestre MR, Palmer C, Toro N, Fineran PC, Jackson SA. 2022. PADLOC: a web server for the
618 identification of antiviral defence systems in microbial genomes. *Nucleic acids research* **50**: W541-
619 W550.
- 620 Pinilla-Redondo R, Russel J, Mayo-Munoz D, Shah SA, Garrett RA, Nesme J, Madsen JS, Fineran PC, Sorensen SJ.
621 2022. CRISPR-Cas systems are widespread accessory elements across bacterial and archaeal plasmids.
622 *Nucleic acids research* **50**: 4315-4328.
- 623 Pinilla-Redondo R, Shehreen S, Marino ND, Fagerlund RD, Brown CM, Sorensen SJ, Fineran PC, Bondy-Denomy
624 J. 2020. Discovery of multiple anti-CRISPRs highlights anti-defense gene clustering in mobile genetic
625 elements. *Nature communications* **11**: 5652.
- 626 Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for large alignments.
627 *PLoS One* **5**: e9490.
- 628 Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin EV. 2014. Genomes in turmoil: quantification of
629 genome dynamics in prokaryote supergenomes. *BMC biology* **12**: 66.
- 630 Puigbo P, Makarova KS, Kristensen DM, Wolf YI, Koonin EV. 2017. Reconstruction of the evolution of microbial
631 defense systems. *BMC evolutionary biology* **17**: 94.
- 632 Pursey E, Dimitriu T, Paganelli FL, Westra ER, van Houte S. 2022. CRISPR-Cas is associated with fewer antibiotic
633 resistance genes in bacterial pathogens. *Philos Trans R Soc Lond B Biol Sci* **377**: 20200464.
- 634 Rocha EPC, Bikard D. 2022. Microbial defenses against mobile genetic elements and viruses: Who defends
635 whom from what? *PLoS Biol* **20**: e3001514.
- 636 Rousset F, Depardieu F, Miele S, Dowding J, Laval AL, Lieberman E, Garry D, Rocha EPC, Bernheim A, Bikard D.
637 2022. Phages and their satellites encode hotspots of antiviral systems. *Cell Host Microbe* **30**: 740-753
638 e745.
- 639 Russel J, Pinilla-Redondo R, Mayo-Munoz D, Shah SA, Sorensen SJ. 2020. CRISPRCasTyper: Automated
640 Identification, Annotation, and Classification of CRISPR-Cas Loci. *CRISPR J* **3**: 462-469.
- 641 Säfken B, Rügamer D, Kneib T, Greven S. 2021. Conditional Model Selection in Mixed-Effects Models with
642 cAIC4. *Journal of Statistical Software* **99**: 1 - 30.
- 643 Shaw LP, Rocha EPC, MacLean RC. 2023. Restriction-modification systems have shaped the evolution and
644 distribution of plasmids across bacteria. *Nucleic acids research* **51**: 6806-6818.
- 645 Shehreen S, Chyou TY, Fineran PC, Brown CM. 2019. Genome-wide correlation analysis suggests different roles
646 of CRISPR-Cas systems in the acquisition of antibiotic resistance genes in diverse species. *Philos Trans*
647 *R Soc Lond B Biol Sci* **374**: 20180384.
- 648 Soucy SM, Huang J, Gogarten JP. 2015. Horizontal gene transfer: building the web of life. *Nature reviews*
649 *Genetics* **16**: 472-482.
- 650 Stamatakis A. 2014. RAXML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.
651 *Bioinformatics* **30**: 1312-1313.
- 652 Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence alignments into the
653 corresponding codon alignments. *Nucleic acids research* **34**: W609-612.
- 654 Tesson F, Herve A, Mordret E, Touchon M, d'Humieres C, Cury J, Bernheim A. 2022. Systematic and quantitative
655 view of the antiviral arsenal of prokaryotes. *Nature communications* **13**: 2561.
- 656 van Houte S, Buckling A, Westra ER. 2016. Evolutionary Ecology of Prokaryotic Immune Mechanisms.
657 *Microbiology and molecular biology reviews : MMBR* **80**: 745-763.
- 658 van Vliet AHM, Charity OJ, Reuter M. 2021. A Campylobacter integrative and conjugative element with a
659 CRISPR-Cas9 system targeting competing plasmids: a history of plasmid warfare? *Microb Genom* **7**.
- 660 Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson P, Weckesser
661 W, Bright J et al. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat*
662 *Methods* **17**: 261-272.
- 663 Watson BNJ, Capria L, Alseth EO, Pons B, Biswas A, Lenzi L, Buckling A, van Houte S, Westra ER, Meaden S.
664 2024. CRISPR-Cas in *Pseudomonas aeruginosa* provides transient population-level immunity against
665 high phage exposures. *The ISME journal* doi:10.1093/ismejo/wrad039.

- 666 Watson BNJ, Staals RHJ, Fineran PC. 2018. CRISPR-Cas-Mediated Phage Resistance Enhances Horizontal Gene
667 Transfer by Transduction. *mBio* **9**.
- 668 Westra ER, Levin BR. 2020. It is unclear how important CRISPR-Cas systems are for protecting natural
669 populations of bacteria against infections by mobile genetic elements. *Proceedings of the National*
670 *Academy of Sciences of the United States of America* **117**: 27777-27785.
- 671 Wheatley RM, MacLean RC. 2021. CRISPR-Cas systems restrict horizontal gene transfer in *Pseudomonas*
672 *aeruginosa*. *The ISME journal* **15**: 1420-1433.
- 673 Wolf YI, Makarova KS, Lobkovsky AE, Koonin EV. 2016. Two fundamentally different classes of microbial genes.
674 *Nature microbiology* **2**: 16208.
- 675 Wu Y, Garushyants SK, van den Hurk A, Aparicio-Maldonado C, Kushwaha SK, King CM, Ou Y, Todeschini TC,
676 Clokie MRJ, Millard AD et al. 2024. Bacterial defense systems exhibit synergistic anti-phage activity.
677 *Cell Host Microbe* **32**: 557-572 e556.
- 678 Yan Y, Zheng J, Zhang X, Yin Y. 2024. dbAPIS: a database of anti-prokaryotic immune system genes. *Nucleic*
679 *acids research* **52**: D419-D425.

680

681

682 **Figure legends**683 **Figure 1: Association between 7 widespread defense systems, total number of genes, and MGE**

684 **abundance.** (a) Number of species displaying positive or negative associations in a phylogenetic
 685 generalized linear mixed effects model (see Methods) according to two different criteria: statistical
 686 significance ($p < 0.05$) and absolute effect size greater than the smallest significant effect size
 687 ($|ES| > SSES$), separately computed for each response variable and defense system. (Note that, due to
 688 cross-species differences in statistical power, the second criterion does not necessarily imply
 689 statistical significance.) The top row (Total) indicates the association with the total number of genes.
 690 The association with MGE was calculated based on marker genes for prophages (Ph), plasmids (Pl),
 691 and transposons (Tr). Abbreviations of functional categories: X (mobilome), L (replication,
 692 recombination and repair), U (intracellular trafficking and secretion), K (transcription), D (cell cycle
 693 control, cell division and chromosome partitioning), and V (defense). (b) Effect sizes, measured as
 694 relative differences in gene and MGE abundances. Each point corresponds to one species. Species
 695 with values beyond the axis limits are collapsed in a single point with size proportional to the number
 696 of species. Vertical lines indicate the median over all the species that show positive or negative
 697 association according to the p-value and SSES criteria.

698 **Figure 2: Taxonomic distribution of species displaying positive and negative associations between**

699 **the presence of defense systems and the number of genes from the mobilome** (based on COG

700 annotations). Taxa discussed in the text are labeled: A, *Acinetobacter*; Pa, *Pseudomonas aeruginosa*;

701 Pv, *Phocaeicola vulgatus*; Pd, *Phocaeicola dorei*; S1, *Streptococcus pyogenes*, *S. dysgalactiae*, *S. equi*,

702 and *S. uberis*; S2, *S. gordonii*, *S. anginosus*, *S. mutans*, and *S. salivarius*; S3, *S. oralis*, *S. intermedius*,

703 and *S. suis*. The classes *Bacilli* and *Clostridia* are the major components of the Genome Taxonomy

704 Database phyla “Bacillota” and “Bacillota_A”. Note that, due to cross-species differences in statistical

705 power, the condition $|ES| > SSES$ does not necessarily imply statistical significance. See figure S3 for

706 a larger version including all species names.

707 **Figure 3: Relative differences in the rates of gene gain and loss between sister clades that do and**
708 **do not harbor defense systems** (DEF^+ and DEF^- , respectively). The boxplots represent the distribution
709 of the DEF^+/DEF^- ratio of gene gain (or loss) rates for each class of MGE, calculated for every pair of
710 sister clades. Values greater (or smaller) than 1 indicate increased (or reduced) gene flux in lineages
711 that contain the defense system. (a) Species that show a negative association between the defense
712 system and the number of marker genes for each class of MGE (PGLMM with smallest significant
713 effect size criterion). (b) Species that show a positive association between the defense system and
714 the number of marker genes for each class of MGE. In the boxplots, the central line indicates the
715 median, the box limits correspond to the 25 and 75 percentiles, and the whiskers extend to the
716 largest and smallest values not classified as outliers. P-values are based on Wilcoxon test with log-
717 ratio = 0 as null hypothesis.

718 **Figure 4: Co-occurrence of defense system gains and losses and MGE gains and losses along the**
719 **phylogeny.** (a) Percentage of DEF gains (and losses) that occur in the same branch as an MGE gain (or
720 loss). (b) Conditional probability of gaining (or losing) a defense system provided that an MGE is also
721 gained (or lost) in the same branch, compared to the conditional probabilities when an MGE is not
722 acquired (or lost) in the same branch. (c) Effect of MGE gain (or loss) on the per branch probability to
723 acquire (or loss) a defense system, measured as a risk ratio. Error bars in (a) and (b) correspond to
724 95% confidence intervals based on the binomial distribution. Error bars in (c) indicate the 95%
725 confidence intervals for the relative risk (Morris and Gardner 1988).

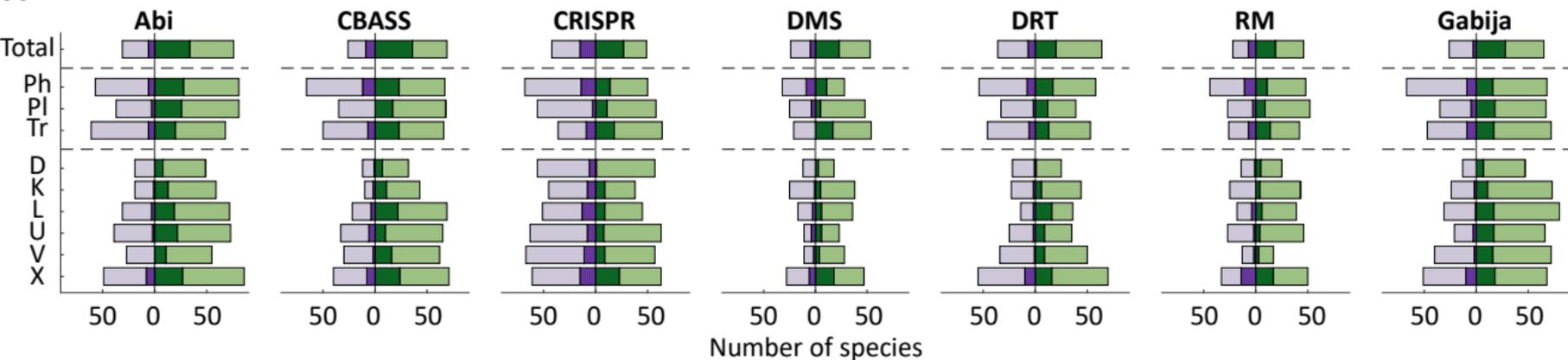
726 **Figure 5: Association between defense systems and gene gain rates depends on the timescale.** (a)
727 Ratio of overall gene gain rates between sister clades that do and do not harbor defense systems
728 (DEF^+ and DEF^- , respectively). Values greater (or smaller) than 1 indicate increased (or reduced) gene
729 flux in lineages that contain the defense system. (b) Density distribution of the DEF^+/DEF^- gain ratios,
730 comparing clades that represent recent acquisitions of the defense system (depth $< 10^{-5}$ substitutions
731 per site in core genes) and clades that have retained the defense system for longer times (depth $> 10^{-5}$

732 ³ substitutions per site). The probability density function (y-axis) is represented in logarithmic scale to
733 facilitate the visualization the tails. The skewness of all distributions is statistically significant (Table
734 S7), with positive values for recent clades and negative values for deeper clades (c) Representative
735 examples of gene gain ratios in shallow and deeper sister groups from the same species. (d)
736 Percentage of genes from different defense systems located within known MGE. Whiskers represent
737 95% confidence intervals based on the binomial distribution.

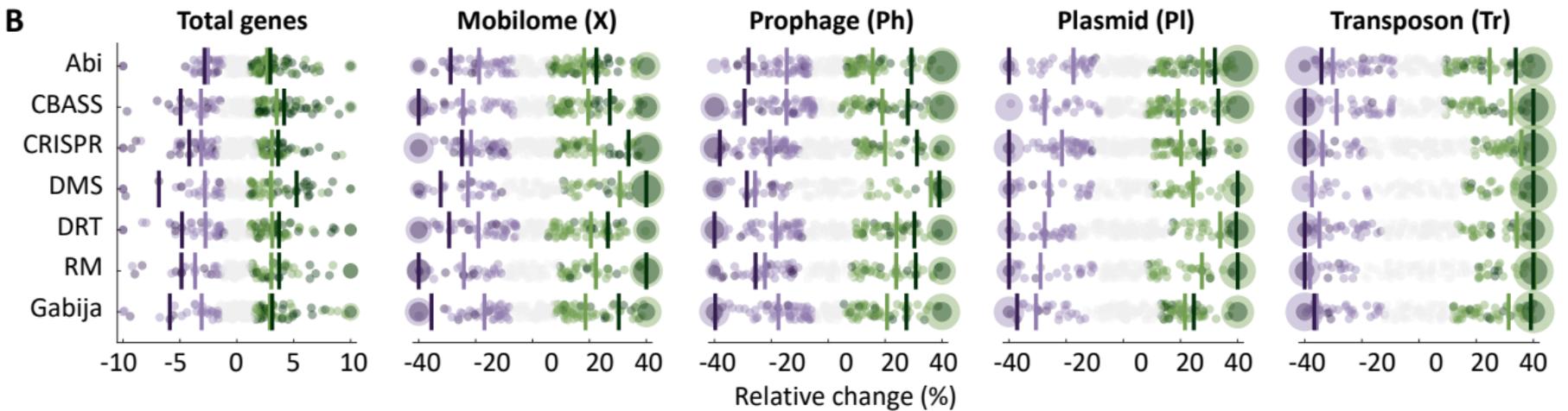
738

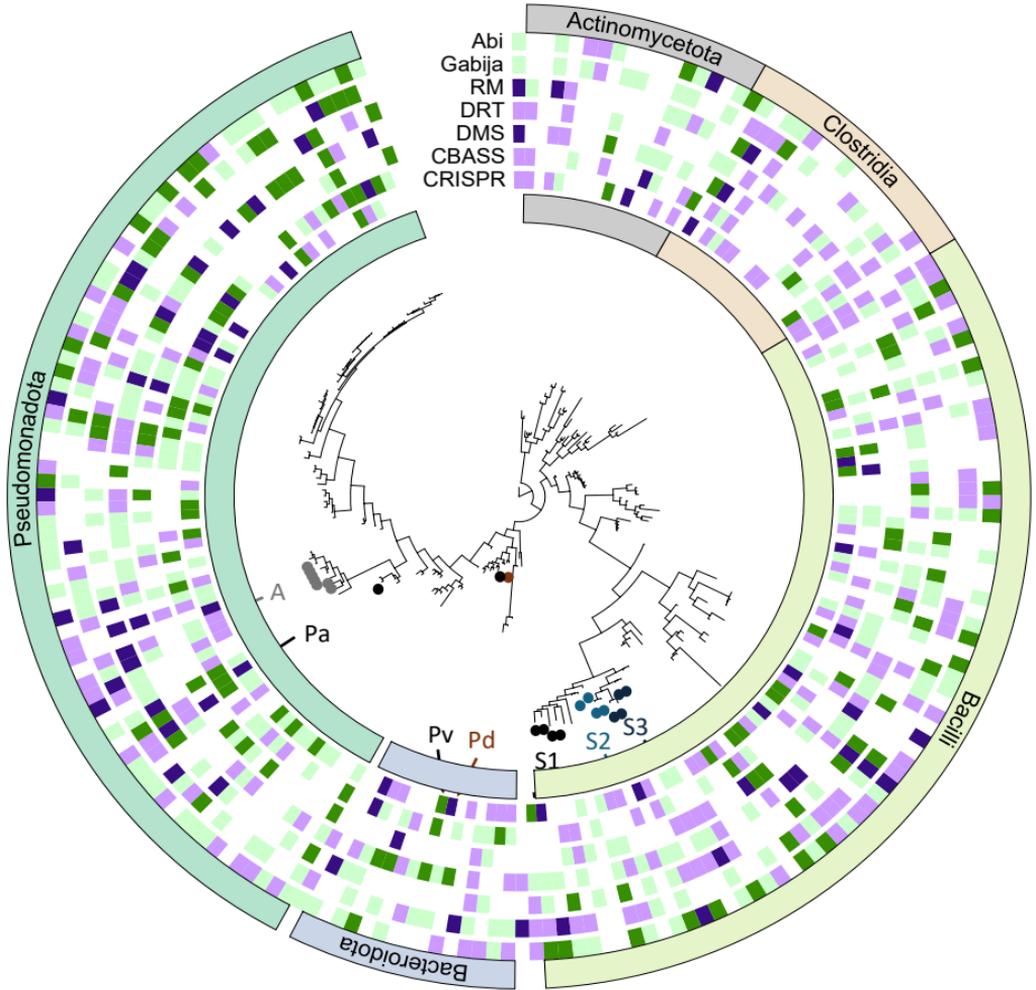
Negative ($|ES| > SSES$)
 Negative ($p < 0.05$)
 Positive ($p < 0.05$)
 Positive ($ES > SSES$)

A



B





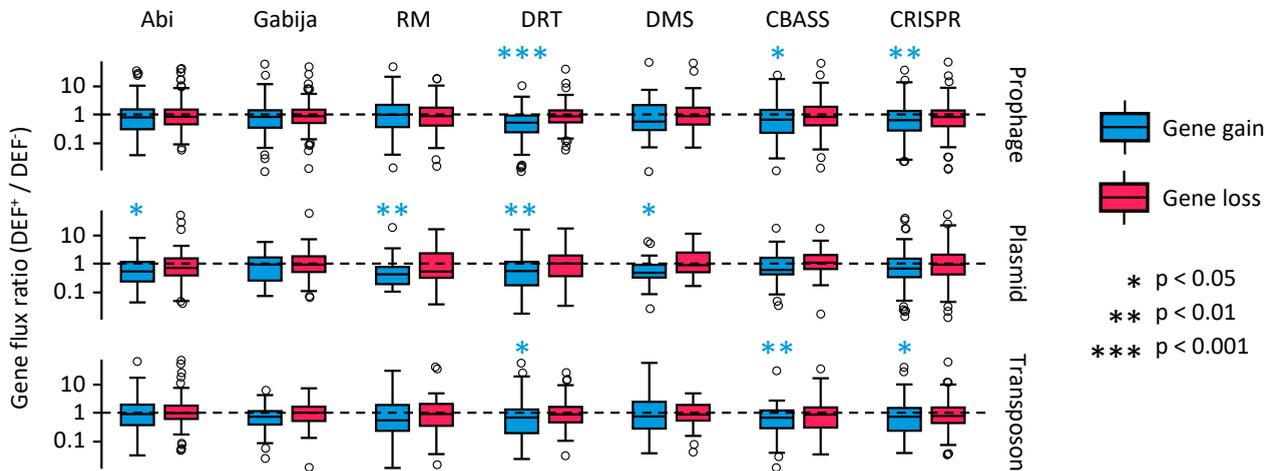
■ Negative ($p < 0.05$)

■ Negative ($|ES| > SSES$)

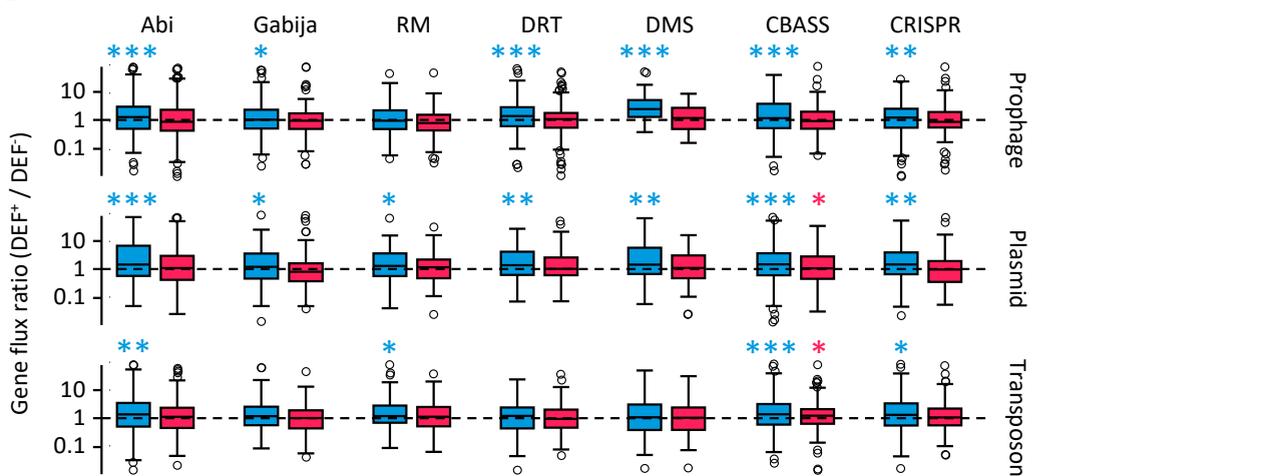
■ Positive ($p < 0.05$)

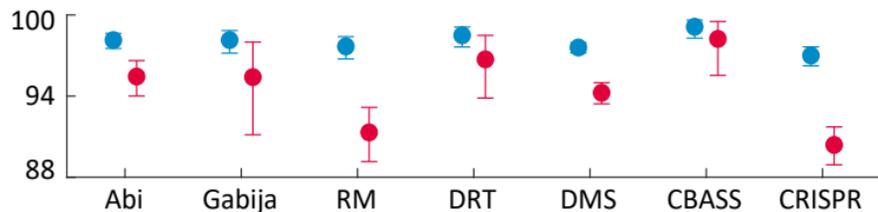
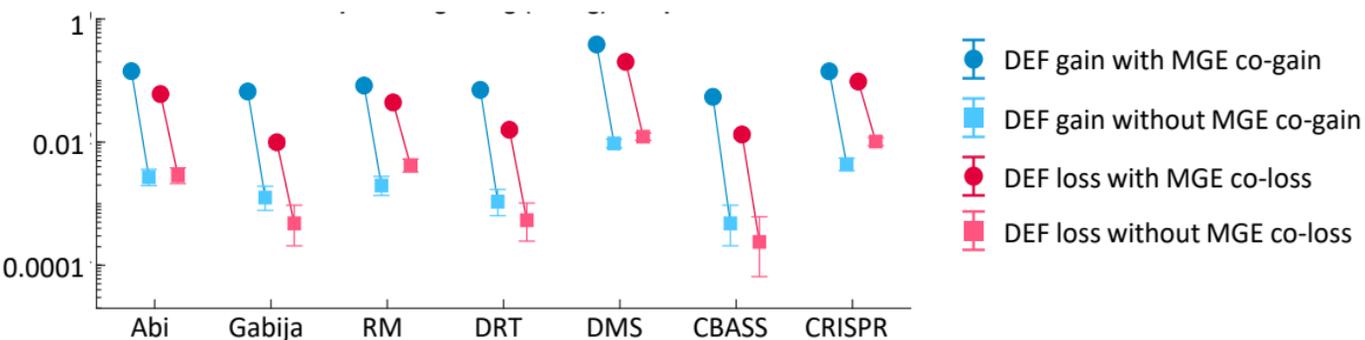
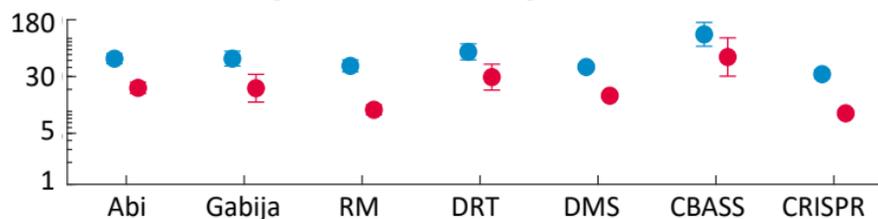
■ Positive ($ES > SSES$)

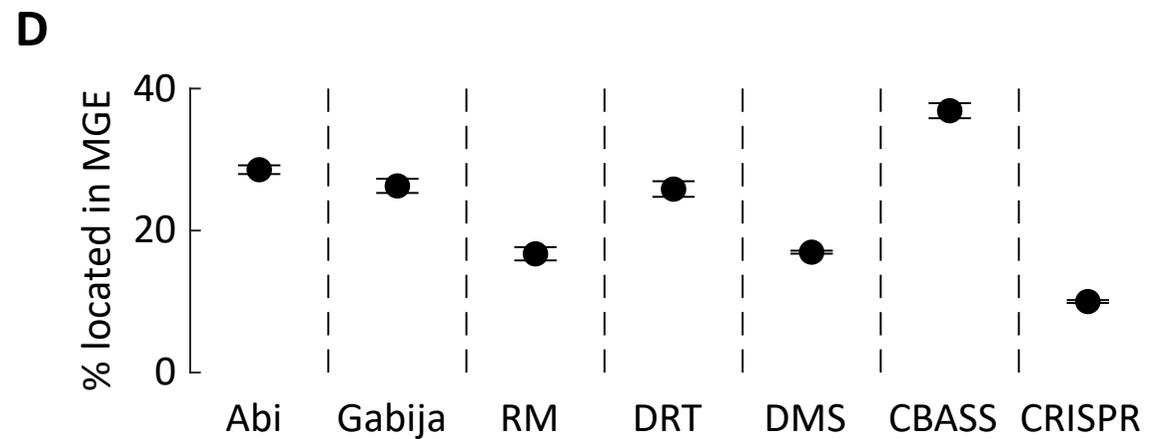
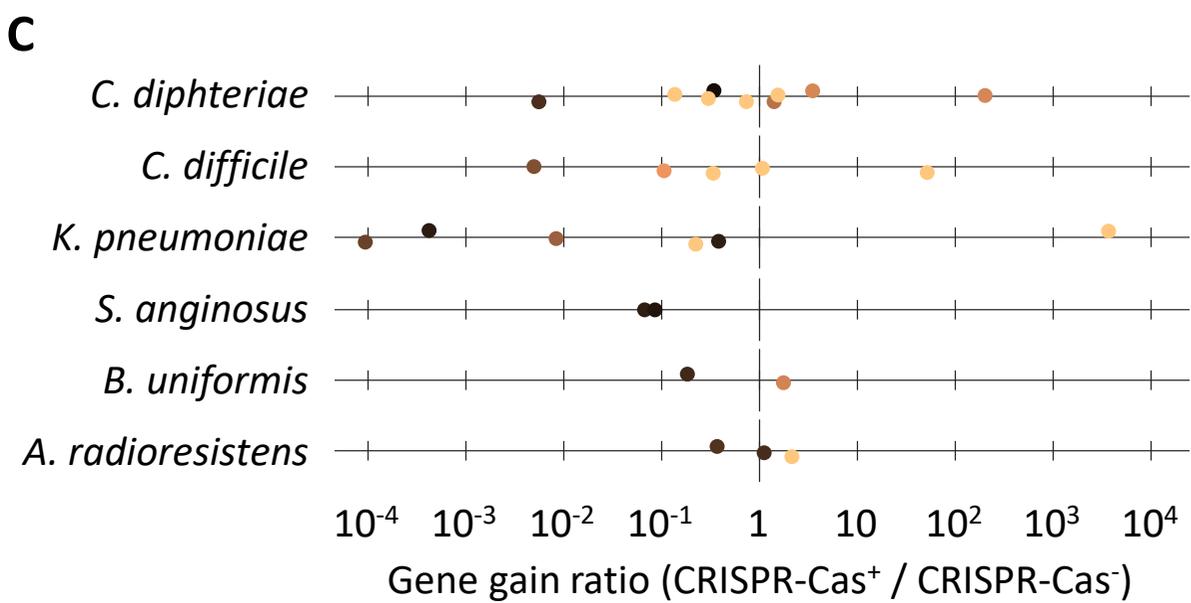
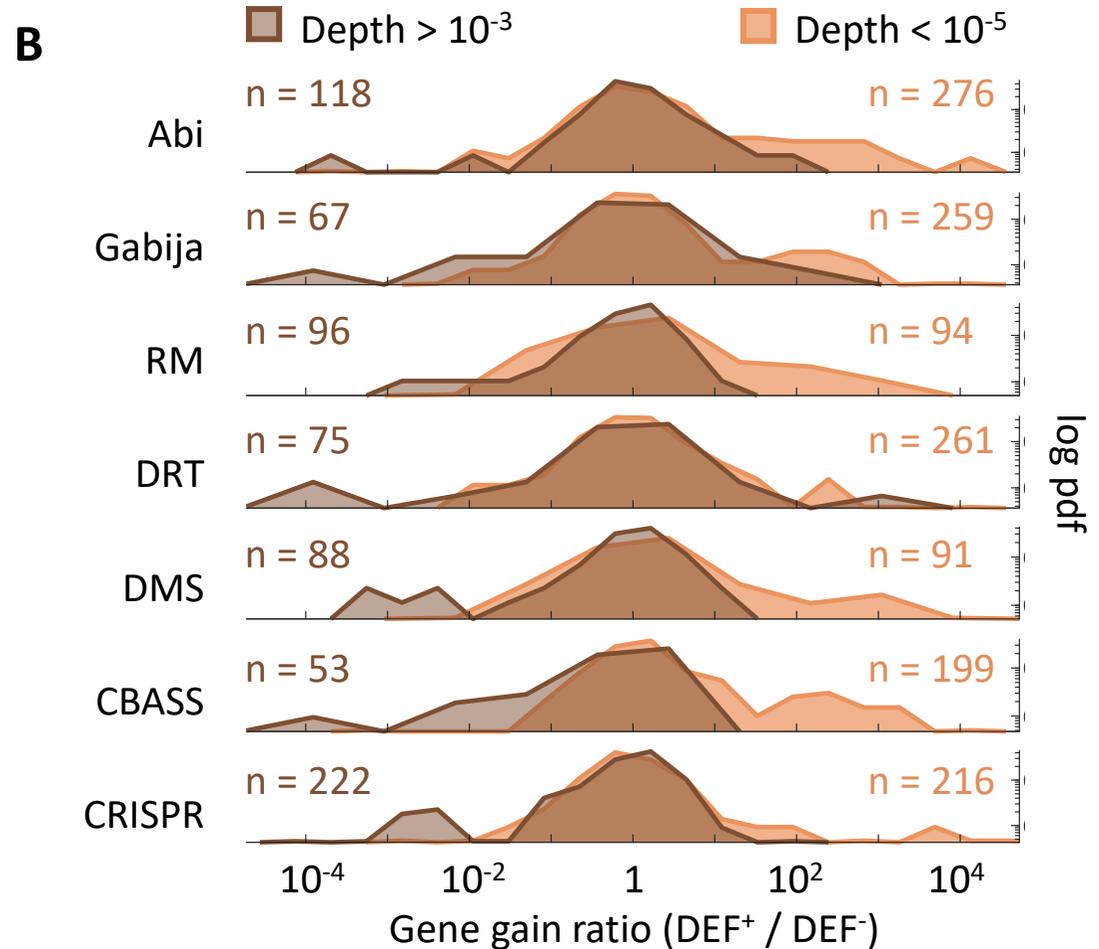
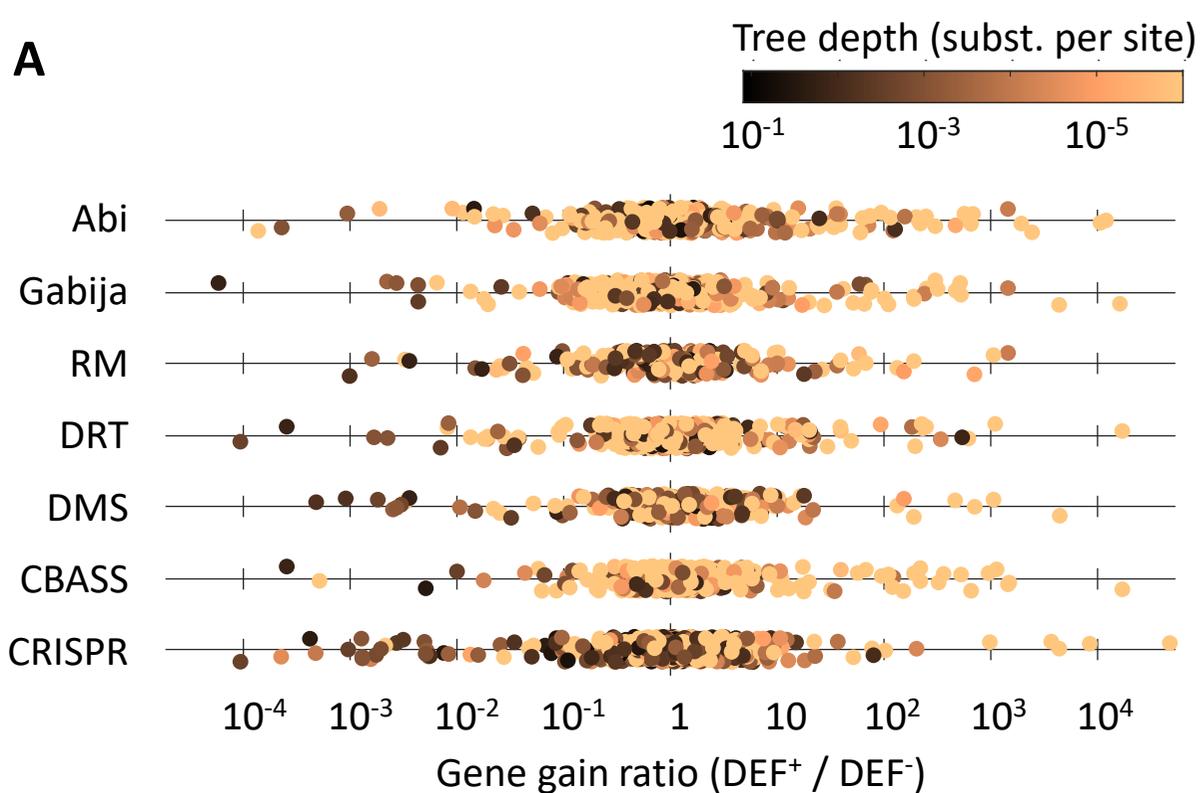
A Species with negative association between defense and MGE abundance

A


B Species with positive association between defense and MGE abundance

B


A DEF gains and losses that co-occur with MGE gain and loss (%)**B** Conditional prob. of gaining or losing DEF per branch**C** Effect of MGE gain and loss on DEF gain and loss (relative risk)





Timescale and genetic linkage explain the variable impact of defense systems on horizontal gene transfer

Yang Liu, Joao Botelho and Jaime Iranzo

Genome Res. published online January 10, 2025

Access the most recent version at doi:[10.1101/gr.279300.124](https://doi.org/10.1101/gr.279300.124)

Supplemental Material <http://genome.cshlp.org/content/suppl/2025/01/30/gr.279300.124.DC1>

P<P Published online January 10, 2025 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This manuscript is Open Access. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International license), as described at <http://creativecommons.org/licenses/by/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
