

1 **Phylogroup-specific variation shapes the clustering of antimicrobial resistance genes and defence**
2 **systems across regions of genome plasticity**

3

4

5 João Botelho^{1,2*}, Leif Tüffers^{2,3}, Janina Fuss⁴, Florian Buchholz², Christian Utpatel^{5,6}, Jens Klockgether⁷,
6 Stefan Niemann^{5,6}, Burkhard Tümmler^{7,8}, Hinrich Schulenburg^{1,2*}

7

8 ¹Antibiotic resistance group, Max-Planck Institute for Evolutionary Biology, Plön, Germany;

9 ²Evolutionary Ecology and Genetics, University of Kiel, Kiel, Germany;

10 ³Department of Infectious Diseases and Microbiology, University of Lübeck, Lübeck, Germany;

11 ⁴Institute of Clinical Molecular Biology, Christian Albrechts University and University Hospital
12 Schleswig-Holstein, Kiel, Germany;

13 ⁵Molecular and Experimental Mycobacteriology, Research Center Borstel, Borstel, Germany;

14 ⁶German Center for Infection Research, Partner Site Hamburg-Lübeck-Borstel-Riems, Borstel, Germany

15 ⁷Clinic for Paediatric Pneumology, Allergology, and Neonatology, Hannover Medical School (MHH),
16 Hannover, Germany;

17 ⁸Biomedical Research in Endstage and Obstructive Lung Disease Hannover (BREATH), German Center
18 for Lung Research, Hannover Medical School, Hannover, Germany.

19

20 To whom correspondence should be addressed. Email: botelho@evolbio.mpg.de. Correspondence may
21 also be addressed to hschulenburg@zoologie.uni-kiel.de.

22

23

24

25

26

27

28

29 **Summary**

30 **Background** *Pseudomonas aeruginosa* is an opportunistic pathogen consisting of three phylogroups
31 (hereafter named A, B, and C) of unevenly distributed size. Here, we assessed phylogroup-specific
32 evolutionary dynamics in a collection of *P. aeruginosa* genomes.

33 **Methods** In this genomic analysis, using phylogenomic and comparative genomic analyses, we generated
34 18 hybrid assemblies from a phylogenetically diverse collection of clinical and environmental *P.*
35 *aeruginosa* isolates, and contextualised this information with 1991 publicly available genomes of the
36 same species. We explored to what extent antimicrobial resistance (AMR) genes, defence systems, and
37 virulence genes vary in their distribution across regions of genome plasticity (RGPs) and “masked”
38 (RGP-free) genomes, and to what extent this variation differs among the phylogroups.

39 **Findings** We found that members of phylogroup B possess larger genomes, contribute a comparatively
40 larger number of pangenome families, and show lower abundance of CRISPR-Cas systems. Furthermore,
41 AMR and defence systems are pervasive in RGPs and integrative and conjugative/mobilizable elements
42 (ICEs/IMEs) from phylogroups A and B, and the abundance of these cargo genes is often significantly
43 correlated. Moreover, inter- and intra-phylogroup interactions occur at the accessory genome level,
44 suggesting frequent recombination events. Finally, we provide here a panel of diverse *P. aeruginosa*
45 strains to be used as reference for functional analyses.

46 **Interpretation** Altogether, our results highlight distinct pangenome characteristics of the *P. aeruginosa*
47 phylogroups, which are possibly influenced by variation in the abundance of CRISPR-Cas systems and
48 that are shaped by the differential distribution of other defence systems and AMR genes.

49 **Funding** German Science Foundation, Max-Planck Society, Leibniz ScienceCampus Evolutionary
50 Medicine of the Lung, BMBF program Medical Infection Genomics, Kiel Life Science Postdoc Award.

51

52

53

54

55

56

57

58

59

60

61 **Research in context**

62 **Evidence before this study**

63 To date, pangenome studies exploring the epidemiology and evolution dynamics of bacterial pathogens
64 have been limited due to the use of gene frequencies across whole species dataset without accounting for
65 biased sampling or the population structure of the genomes in the dataset. We searched PubMed without
66 language restrictions for articles published before September 1, 2021, that investigated the phylogroup-
67 specific evolutionary dynamics across bacterial species. In this literature search we used the search terms
68 “pangenome” and “phylogroup” or “uneven”, which returned 14 results. Of these, only one study used a
69 population structure-aware approach to explore pangenome dynamics in a bacterial species consisting of
70 multiple phylogroups with unevenly distributed members.

71 **Added value of this study**

72 To our knowledge, this study is the first to assess phylogroup-specific evolutionary dynamics in a
73 collection of genomes belonging to the nosocomial pathogen *P. aeruginosa*. Using a refined approach
74 that challenges traditional pangenome analyses, we found specific signatures for each of the three
75 phylogroups, and we demonstrate that members of phylogroup B contribute a comparatively larger
76 number of pangenome families, have larger genomes, and have a lower prevalence of CRISPR-Cas
77 systems. Additionally, we observed that antibiotic resistance and defence systems are pervasive in regions
78 of genome plasticity and integrative and conjugative/mobilizable elements from phylogroups A and B,
79 and that antibiotic resistance and defence systems are often significantly correlated in these mobile
80 genetic elements.

81 **Implications of all the available evidence**

82 These results indicate that biases inherent to traditional pangenome approaches can obscure the real
83 distribution of important cargo genes in a bacterial species with a complex population structure.
84 Furthermore, our findings pave the way to new pangenome approaches that are currently under-explored
85 in comparative genomics and, crucially, shed a new light on the role that integrative and
86 conjugative/mobilizable elements may play in protecting the host against foreign DNA.

87

88

89

90

91

92

93

94 Introduction

95 *Pseudomonas aeruginosa* is a ubiquitous metabolically versatile γ -proteobacterium. This Gram-negative
96 bacterium is also an opportunistic human pathogen commonly linked to life-threatening acute and chronic
97 infections ¹. It belongs to the ESKAPE pathogens collection ², highlighting its major contribution to
98 nosocomial infections across the globe and its ability to “escape” antimicrobial therapy because of the
99 widespread evolution of antimicrobial resistance (AMR) ³. This species is also often found to be multi- as
100 well as extensively drug resistant (MDR and XDR, respectively) ⁴, making it difficult and in some cases
101 even impossible to treat. For this reason, *P. aeruginosa* is placed by the World Health Organization
102 (WHO) in the top priority group of most critical human pathogens, for which new treatment options are
103 urgently required ⁵. These efforts rely on an in-depth understanding of the species biology and its
104 evolutionary potential, which may be improved through a functional analysis of whole genome
105 sequencing data.

106 The combined pool of genes belonging to the same bacterial species is commonly referred to as the
107 pangenome. Frequently, only a small proportion of these genes is shared by all species members (the core
108 genome). On the contrary, a substantial proportion of the total pool of genes is heterogeneously
109 distributed across the members (the accessory genome). Following Koonin and Wolf ⁶, the pangenome
110 can be divided into 3 categories: i) the persistent or softcore genome, for gene families present in the
111 majority of the genomes; ii) the shell genome, for those present at intermediate frequencies and that are
112 gained and lost rather slowly; iii) the cloud genome, for gene families present at low frequency in all
113 genomes and that are rapidly gained and lost ⁷. Clusters of genes that are part of the accessory genome
114 (i.e, the shell and cloud genome) are often located in so-called regions of genome plasticity (RGPs),
115 genomic loci apparently prone to insertion of foreign DNA. By harbouring divergent accessory DNA in
116 different strains, these loci can represent highly variable genomic regions. The shell and cloud genomes
117 are also characterized by mobile genetic elements (MGEs) that are capable of being laterally transferred
118 between bacterial cells, including plasmids, integrative and conjugative/mobilizable elements
119 (ICEs/IMEs), and prophages ^{8,9}. These MGEs can mediate the shuffling of cargo genes that may provide a
120 selective advantage to the recipient cell, such as resistance to antibiotics, increased pathogenicity, and
121 defence systems against foreign DNA ¹⁰⁻¹².

122 Most pangenome studies described to date have characterized gene frequencies across the whole species
123 dataset without accounting for biased sampling or the population structure of the genomes in the dataset.
124 This is particularly relevant in species consisting of multiple phylogroups with unevenly distributed
125 members. As recently reported for *Escherichia coli* ¹³, genes classified as part of the accessory genome
126 using traditional pangenome approaches are in fact core to specific phylogroups. Since *P. aeruginosa* is
127 composed of three different-sized phylogroups (hereafter referred to as phylogroups A, B, and C as per
128 the nomenclature proposed by Ozer *et al* ¹⁴; see also results), characterized by high intraspecies functional
129 variability ^{15,16}, it is likely that evolution in these phylogroups is driven by specific sets of genes found in
130 the majority of members within the groups, but not across groups.

131 The aim of the current study is to enhance our understanding of the pangenome of the human pathogen *P.*
132 *aeruginosa* by specifically assessing phylogroup-specific characteristics and genome dynamics, including
133 data from more than 2000 genomes. We explore to what extent particular groups of cargo genes, such as
134 those encoding AMR, virulence, and defence systems, vary in their distribution across RGP and
135 “masked” (RGP-free) genomes, and to what extent this variation differs among the phylogroups. Our data
136 set includes new full genome sequences of a representative set of *P. aeruginosa* strains, the ‘major *P.*
137 *aeruginosa* clone type’ (mPact) strain panel. This set of strains was previously isolated by the Tümmler
138 lab (Hanover, Germany) from both clinical and environmental samples¹⁷. This mPact panel encompasses
139 the most common clone types in the contemporary population^{18–20} and provides a manageable, focused
140 resource for in-depth functional analyses.

141

142 **Methods**

143 **Sequencing and hybrid assembly of the mPact strain panel**

144 Genomic DNA from 18 strains of the mPact panel¹⁷ were extracted using the Macherey-Nagel
145 NucleoSpin Tissue kit, according to the standard bacteria support protocol from the manufacturer. We
146 used Nanodrop 1000 for DNA quantification and quality control (260/280 and 260/230 ratios), followed
147 by measurements in Qubit for a more precise quantification. The Agilent TapeStation and the
148 FragmentAnalyzer Genomic DNA 50KB kit served to control fragment size. Sequencing libraries were
149 prepared with the Illumina Nextera DNA flex and Pacific Bioscience (PacBio) SMRTbell express
150 template prep kit 2.0. Libraries were sequenced on the Illumina MiSeq at 2x300bp or the PacBio Sequel
151 II, respectively. Illumina reads were verified for quality using FastQC v0.11.9²¹ and trimmed with Trim
152 Galore v0.6.6²², using the paired-end mode with default parameters and a quality Phred score cutoff of
153 10. Both datasets were then combined using the Unicycler v0.4.8 assembly pipeline²³. We used the
154 default normal mode in Unicycler to build the assembly graphs of most strains, except of the mPact
155 strains H02, H14, H15, H18, and H19, where we used the bold mode. The assemblies were visually
156 inspected using the assembly graph tool Bandage v0.8.1²⁴.

157

158 **Bacterial collection**

159 We downloaded a total of 5468 *P. aeruginosa* genomes from RefSeq’s NCBI database using PanACoTA
160 v1.2.0²⁵. After quality control to remove low-quality assemblies, 2704 were kept and 2764 genomes with
161 more than 100 contigs were discarded (**Table S1**). Next, 713 genomes were discarded by the distance
162 filtering step, using minimum ($1e-4$) and maximum (0.05) mash distance cut-offs to remove duplicates
163 and misclassified assemblies at the species level²⁶, respectively. This resulted in 1991 publicly available
164 genomes. The 18 genomes sequenced in this study from mPact panel¹⁷ passed both filtering steps,
165 resulting in a pruned collection of 2009 genomes in total. Multi-locus sequence typing (MLST) profiles
166 were determined with mlst v2.19.0 (<https://github.com/tseemann/mlst>).

167

168 **Pangenome and phylogenomics**

169 The average nucleotide identity (ANI) between the 2009 genomes was calculated with fastANI v1.33 ²⁶.
170 We used the genome sequences to generate a pangenome with the panrgp subcommand of PPanGGOLiN
171 v1.1.136 ^{27,28}. We built a softcore-genome alignment (threshold 95%), followed by inference of a
172 maximum likelihood tree with the General Time Reversible model of nucleotide substitution in IQ-TREE
173 v2.1.2 ²⁹. To detect recombination events in our collection and account for them in phylogenetic
174 reconstruction, we used ClonalFrameML v1.12 ³⁰. Phylogenetic trees were plotted in iTOL v6
175 (<https://itol.embl.de/>) ³¹ and explored to cluster genomes according to the phylogroup. Due to the sample
176 size difference, we subsequently focused the analysis on each phylogroup separately. Pangenome analysis
177 was performed for each phylogroup, using the panrgp subcommand from PPanGGOLiN. Core and
178 accessory genes were classified across genomes from different phylogroups with a publicly available R
179 script (<https://github.com/ghoresh11/twilight>) ¹³. We used the gene presence/absence output from the
180 whole collection's pangenome and the grouping of our genomes according to the phylogroup.

181

182 **Identification of RGPs and ICEs/IMEs**

183 To mask all the genomes, we used the RGPs coordinates determined by panrgp for each individual
184 genome as input in bedtools maskfasta v2.30.0 ³². We extracted the RGP nucleotide sequences with the
185 help of bedtools getfasta. All genomes were annotated with prokka v1.4.6 ³³. To look for ICEs/IMEs on
186 complete genomes, we used the genbank files created by prokka as input in the standalone-version of
187 ICEfinder ³⁴.

188

189 **Identification of ICEs and functional categories**

190 We retrieved the annotated proteins for the RGPs and masked genomes across the three phylogroups. We
191 clustered each of the six groups of proteins with MMseqs2 v13.45111 ³⁵ and an identity cut-off of 80%.
192 These clustered proteins were scanned for functional categories in eggNOG-mapper v2 ³⁶, using the built-
193 in database for clusters of orthologous groups ³⁷. We calculated the relative frequency of these categories
194 by dividing the absolute counts for each category by the total number of clustered proteins found in each
195 of the six groups. CRISPR-Cas systems were identified with the help of CRISPRCasTyper v1.2.3 ³⁸.
196 AMRfinder v3.10.18 ³⁹ served to locate AMR genes and resistance-associated point mutations.
197 Virulence genes were characterized with the pre-downloaded database from VFDB ⁴⁰ (updated on the 12-
198 05-2021 and including 3867 virulence factors) in abricate v1.0.1 (<https://github.com/tseemann/abricate>).
199 Finally, we searched for defence systems using the protein sequences generated by prokka as input in
200 defense-finder v0.0.11 ⁴¹, a tool developed to identify known defence systems in prokaryotic genomes, for
201 which at least one experimental evidence of the defence function is available.

202

203 **Network-based analysis of RGPs and ICEs/IMEs**

204 As a first step, we calculated the Jaccard Index between the RGPs with the help of BinDash v0.2.1⁴² with
205 k -mer size equal to 21 bp. In detail, we used the sketch subcommand to reduce multiple sequences into
206 one sketch, followed by the dist subcommand, to estimate distance (and relevant statistics) between RGPs
207 in query sketch and RGPs in target-sketch. The Jaccard Index between ICEs/IMEs was similarly obtained
208 with BinDash. We used the mean() function in R to calculate the arithmetic mean of the Jaccard Index.
209 Only Jaccard Index values equal to or above the mean were considered, and the mutation distances served
210 as edge attributes to plot the networks with Cytoscape v3.9.1 under the prefuse force directed layout
211 (<https://cytoscape.org/>). Based on the Analyzer function in Cytoscape, we computed a comprehensive set
212 of topological parameters, such as the clustering coefficient, the network density, the centralization, and
213 the heterogeneity. Clusters in our networks were identified with the AutAnnotate and clusterMaker apps
214 available in Cytoscape, using the connected components as the clustering algorithm.

215

216 **Statistical analysis**

217 The correlation matrix was ordered using the hclust function in R. Statistical comparison of the variation
218 between groups was always based on non-parametric tests, thereby taking into account that the compared
219 groups varied in data distributions (e.g., at least one group with a skewed distribution) and/or showed
220 unequal variances. Moreover, as non-parametric tests are usually considered to be conservative, the thus
221 identified significant test results should indicate trustworthy differences between groups. In particular, the
222 three phylogroups (e.g., genome size, GC content) were generally compared using the Kruskal-Wallis
223 test. The unpaired two-sample Wilcoxon test was used in multiple comparisons between two independent
224 groups of samples (RGPs vs. masked genomes, CRISPR-Cas positive vs negative genomes). In both tests,
225 p -values were adjusted using the Holm–Bonferroni method. Values above 0.05 were considered as non-
226 significant (ns). We used the following convention for symbols indicating statistical significance: * for p
227 ≤ 0.05 , ** for $p \leq 0.01$, *** for $p \leq 0.001$, and **** for $p \leq 0.0001$.

228

229 **Role of funding source**

230 The funders (Max-Planck Society, German Science Foundation, German ministry for education and
231 research, Kiel university) had no role in the study design, data collection, data analysis, data
232 interpretation, or writing of the report. The corresponding authors had full access to all the data and final
233 responsibility for the decision to submit for publication. None of the contributing authors were precluded
234 from accessing data of this study.

235

236

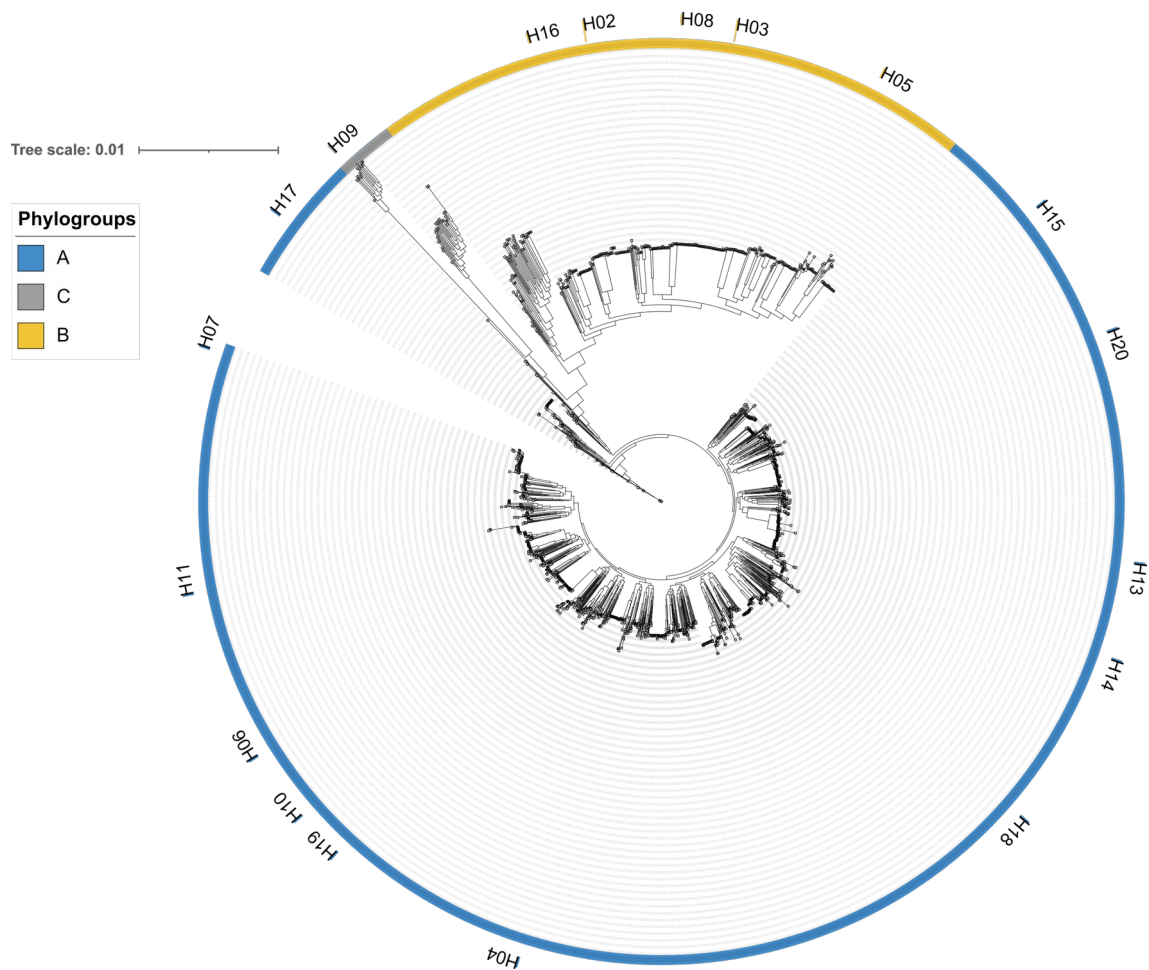
237

238

239 **Results**

240 **The *P. aeruginosa* phylogeny is composed of three phylogroups**

241 Our phylogenomic characterization was based on 2009 assembled *P. aeruginosa* genomes belonging to
242 519 MLST profiles and including 1991 publicly available genomes (following quality control and
243 distance filtering, **Table S2**) and additionally 18 genomes for the mPact strain panel (**Table S3**)¹⁷.
244 Analysis of the ANI values (**Figure S1**) and the softcore-genome alignment of these genomes identified
245 three phylogroups, as previously reported^{14,43} (**Figure 1**). The two major reference isolates are part of the
246 larger phylogroups: PAO1⁴⁴ is part of phylogroup A (n=1531), while the PA14 strain falls into
247 phylogroup B (n=435). Phylogroup C includes a substantially smaller number of members (n=43) (**Table**
248 **S2**). Members of the phylogroup C were recently subdivided into either 2¹⁴ or 3 clusters, including the
249 distantly related PA7 cluster⁴³. In this work, however, the PA7 cluster was excluded, and we focused our
250 analysis on only the remainder of phylogroup C, since genomes from the PA7 cluster were too distantly
251 related to the other genomes. In fact, the PA7 strain was first described as a taxonomic outlier of this
252 species⁴⁵, and genomes belonging to this cluster were recently proposed to belong to a new *Pseudomonas*
253 species⁴⁶. To test the impact of recombination on the softcore-genome alignment, we used
254 ClonalFrameML to reconstruct the phylogenomic tree with corrected branch lengths. The segregation of
255 *P. aeruginosa* into three phylogroups was maintained, resulting in a tree with decreased branch lengths
256 and with identical number of members assigned to each phylogroup (**Figure S2**). Genomes from the
257 mPact panel sequenced in this study were widely distributed across the *P. aeruginosa* phylogeny, with 12
258 strains in phylogroup A, 5 in phylogroup B, and 1 in phylogroup C (**Figure 1** and **Table S3**). Our results
259 show that *P. aeruginosa* consists of three asymmetrical phylogroups and that the segregation of the 2009
260 genomes into phylogenetically distinct groups is not an artefact of recombination.



261

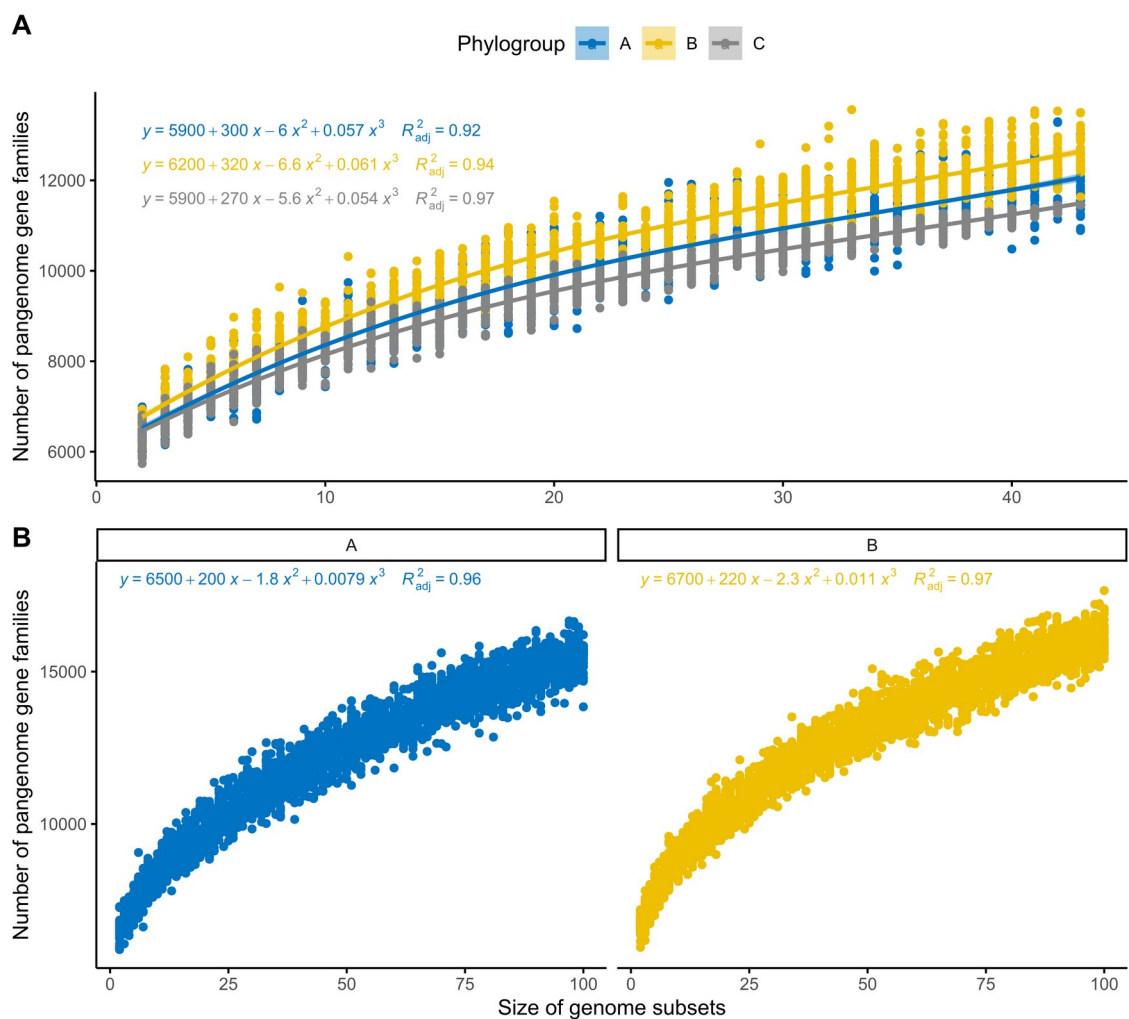
262 **Figure 1.** Maximum-likelihood tree of the softcore-genome alignment of all *P. aeruginosa* isolates used
263 in this study (n=2009). The scale bar represents the genetic distance. Arcs in blue represent phylogroup
264 A, yellow B, and grey C. The phylogenetic placement of the major *P. aeruginosa* clone type (mPact)
265 strain panel, sequenced in this study, are highlighted in the tree, with the strain name (the “H” before each
266 number stands for Hanover, referring to the location of the Tuemmler lab and the study that first
267 described this collection ¹⁷) next to strips coloured according to the phylogroup.

268

269 **Phylogroup B contributes comparatively more gene families to the pangenome than the other two**
270 **phylogroups**

271 We next built a pangenome for the whole species, and separate pangenomes for each of the three
272 phylogroups. This latter approach is important to take phylogenetic subdivisions of the species into
273 account, which is additionally critical because the three phylogroups in our collection have substantially
274 different sample sizes. We observed that the number of persistent gene families in the larger phylogroups
275 A and B were similar to those found in the whole species, while the phylogroup C contained a
276 substantially smaller number of persistent gene families (**Table S4**).

277 The pangenome of bacterial species is usually classified in two types: open pangenomes and closed ones
278 ⁴⁷. Since *P. aeruginosa* is an example of a bacterial species with open pangenome ¹⁴, i.e., the sequencing
279 of new genomes will increase pangenome size, we explored the contribution of each phylogroup to the
280 pangenome. To ensure comparability among the three phylogroups in our first analysis, we randomly
281 drew 43 genomes from each phylogroup (thus, including the total sample size of the smallest phylogroup
282 C), and observed that there is more diversity in the accessory genes of phylogroup B as to the functions
283 contributed by the acquired genes (**Figure 2A** and **Table S5**). In our second analysis, we focused only on
284 the two larger phylogroups A and B, for which we randomly drew in each case 100 genomes and found
285 the trend unchanged (**Figure 2B** and **Table S6**).



286

287 **Figure 2.** Rarefaction curves of the pangenome gene families for each phylogroup. All curves were
288 inferred using polynomial regression lines. Curves in blue represent phylogroup A, yellow B, and grey C.
289 **A)** The curves were generated by randomly re-sampling 43 genomes from each phylogroup several times
290 and then plotting the average number of pangenome families found on each genome. **B)** Rarefaction
291 curves were plotted with 100 random genomes from phylogroups A and B.

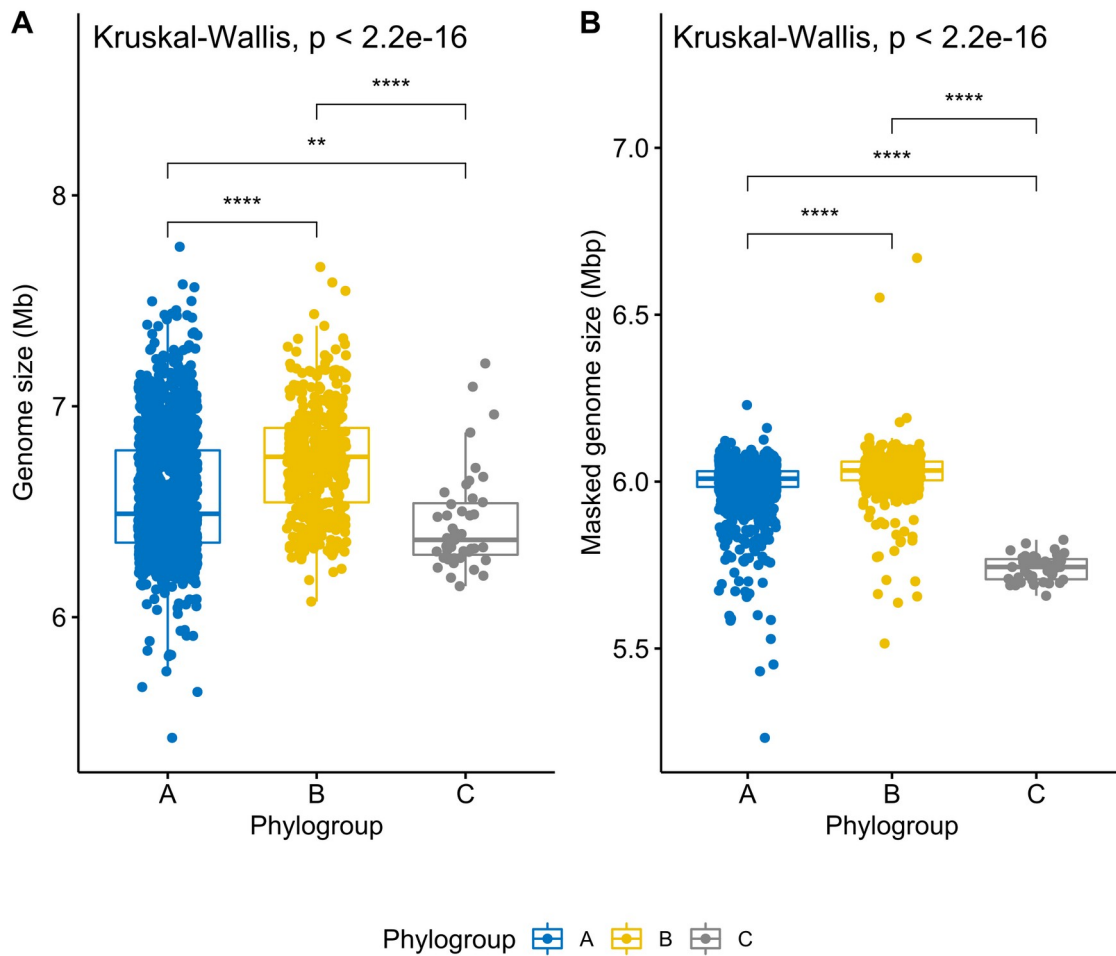
292 We then explored if specific gene families were pervasive across single or multiple phylogroups. We
293 found 14 phylogroup-specific softcore gene families in phylogroup C, and one gene family each was

294 exclusively found in the softcore genomes of phylogroups A and B, respectively (**Figure S3** and **Table**
295 **S7**). Most gene families uniquely found on the softcore genome of phylogroup C were part of the Xcp
296 type-II secretion system (T2SS), which is one of two complete and functionally distinct T2SS present in
297 this species (**Table S8**). The Xcp system is encoded in a cluster containing 11 genes (*xcpP-Z*), as well as
298 an additional *xcpA/pilD* gene found elsewhere in the genome⁴⁸. These genes were also found in the
299 majority of the genomes from phylogroups A and B (**Table S9**), but the encoded proteins were too
300 distantly related to those from phylogroup C. A similar pattern was observed for the two gene families
301 indicated exclusively for either phylogroup A or B, for which we also found distantly related orthologues
302 in phylogroup C. Altogether, these results highlight that phylogroup B differs from the other two by
303 contributing a comparatively larger number of gene families to the pangenome, possibly suggesting that
304 phylogroup B members have larger genomes.

305

306 **Phylogroup B genomes are significantly larger and most carry no CRISPR-Cas systems**

307 A comparison of genome lengths revealed significantly larger genome sizes for phylogroup B than the
308 other two phylogroups (**Figure 3A**, p-value < 2.2e-16). We then extracted the RGPs from each
309 phylogroup, and found a total of 57901 RGPs across the three phylogroups. The RGPs from phylogroup
310 B were significantly larger than that of phylogroup A (**Figure S4**), thus at least contributing to the overall
311 size difference. Nevertheless, after removing the RGPs, the resulting “masked” genomes from phylogroup
312 B were still significantly larger than those from the other two phylogroups (**Figure 3B**, p-value < 2.2e-
313 16). Additionally, we found that genomes from phylogroup B are still significantly larger than those from
314 the other two phylogroups, even if phylogroup sample sizes were adjusted to sample size of the smallest
315 group, phylogroup C (with 43 genomes; **Figure S5**, p-value 3.2e-07). These results point to a potentially
316 higher number of genes conserved across genomes from phylogroup B. Still, the difference in genome
317 size between phylogroups A and B is mainly explained by differences in accessory genome size (**Figure**
318 **S4**). Masked genomes from phylogroup C are significantly smaller than genomes from the other two
319 phylogroups, which is consistent with the smaller number of persistent gene families identified in this
320 phylogroup (**Table S4**). We further explored variation in GC content and observed that the GC content
321 from phylogroup B genomes is significantly lower than those from other phylogroups (**Figure S6**, p-
322 value < 2.2e-16).

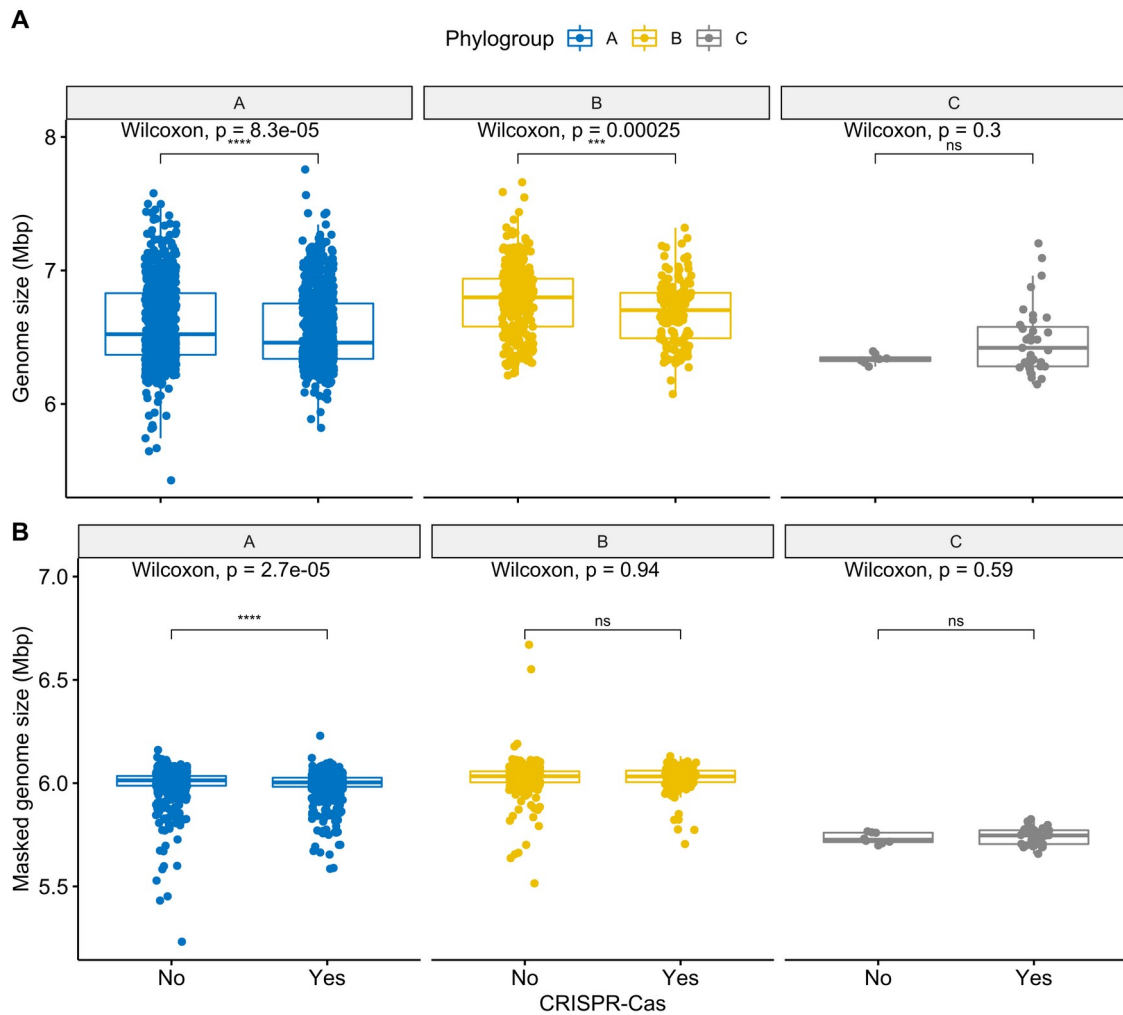


323

324 **Figure 3.** Boxplots representing the variation in genome size (A) and masked genome size (B) across the
325 three phylogroups. Values above 0.05 were considered as non-significant (ns). Stars indicate significance
326 level: * $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq 0.001$, and **** $p \leq 0.0001$. Boxplots in blue represent
327 phylogroup A, yellow B, and grey C.

328 We next assessed whether presence of the defence CRISPR-Cas system is associated with genome size
329 variation. Since CRISPR-Cas systems are important to defend bacteria against foreign DNA^{12,49}, we
330 expected that genomes carrying these systems would be smaller, while those devoid of these systems
331 would accumulate mobile elements and hence be larger. We subdivided genomes from each of the three
332 phylogroups into two groups depending on whether they contain or lack CRISPR-Cas systems,
333 respectively (CRISPR-Cas^{pos}, CRISPR-Cas^{neg}). We indeed found that genomes with CRISPR-Cas systems
334 are significantly smaller than those without (Figure 4A, p-values 8.3e-05 and 0.00025 for the phylogroup
335 A and B comparisons, respectively), supporting the hypothesis that CRISPR-Cas systems can constrain
336 horizontal gene transfer in *P. aeruginosa*⁵⁰⁻⁵². While the number of CRISPR-Cas^{pos} and CRISPR-Cas^{neg}
337 genomes in phylogroups A and C is evenly distributed, phylogroup B genomes without CRISPR-Cas
338 (n=279) were nearly two times more prevalent than those that carried these systems (n=156, Table S2).
339 Interestingly, masked genome size of CRISPR-Cas^{pos} and CRISPR-Cas^{neg} phylogroup B isolates was no
340 longer significantly different from one another (Figure 4B). In line with this finding, we observed that the

341 cumulative size of all RGP was higher in genomes without CRISPR-Cas systems across phylogroups A
342 and B (**Figure S7**). The absence of these defence systems in most genomes from phylogroup B may help
343 to explain the observed larger size.



344

345 **Figure 4.** Boxplots representing the variation in genome size (**A**) and masked genome size (**B**) across
346 pairs of conspecific genomes from the same phylogroup with and without CRISPR-Cas systems. Values
347 above 0.05 were considered as non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq$
348 0.01 , *** $p \leq 0.001$, and **** $p \leq 0.0001$. Boxplots in blue represent phylogroup A, yellow B, and
349 grey C.

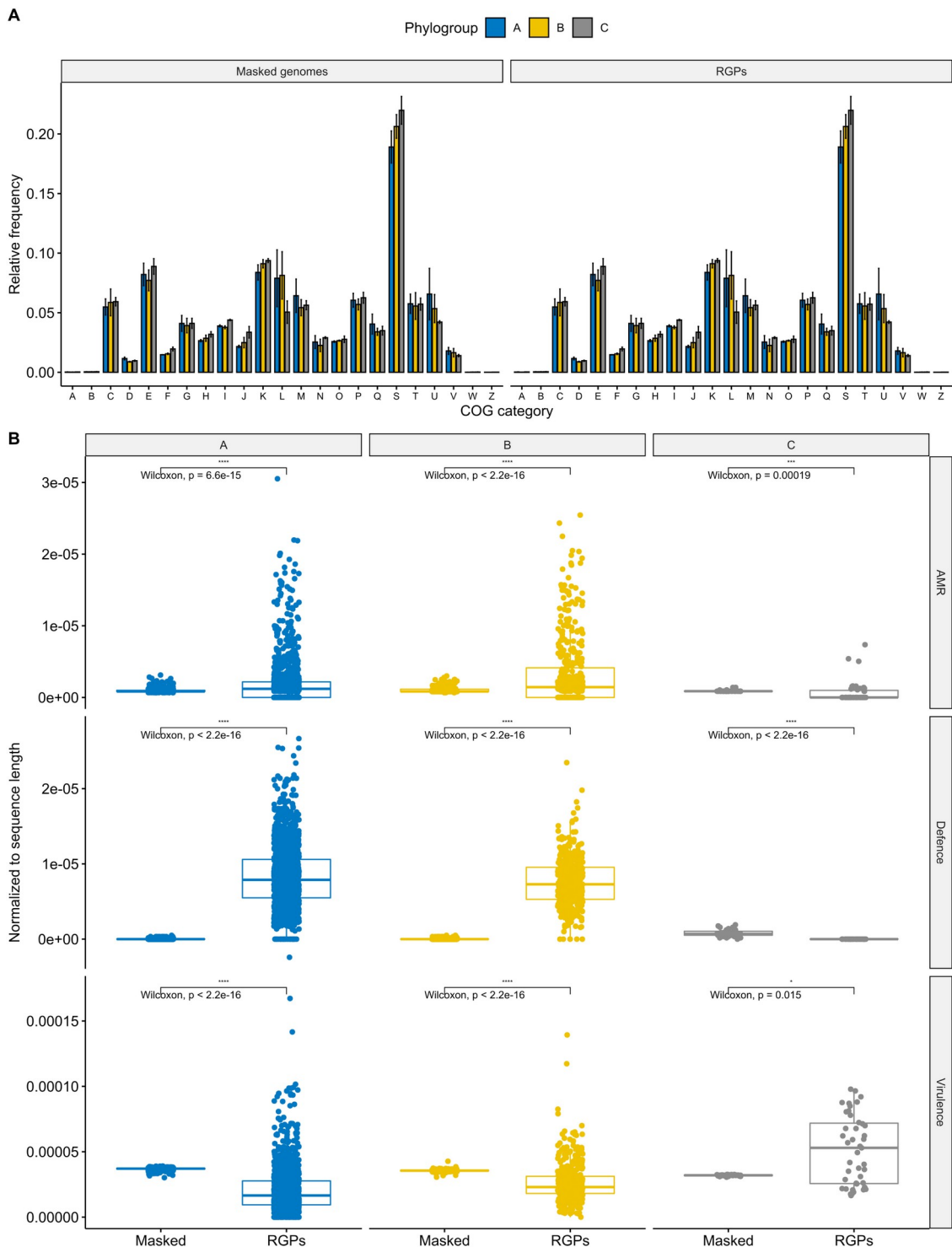
350 We observed a wider diversity of CRISPR-Cas systems in genomes from phylogroup A, including I-C, I-
351 E, I-F, IV-A1, and IV-A2 (**Figure S8** and **Table S10**). These CRISPR-Cas subtypes were all found in
352 genomes from phylogroup B, with the exception of the IV-A2. Curiously, only subtypes I-E and I-F were
353 present in phylogroup C. Type IV CRISPR-Cas systems were found almost exclusively on plasmids, and
354 recent work revealed that they participate in plasmid-plasmid warfare^{12,53}. The type I-C CRISPR-Cas
355 subtype is typically encoded on ICEs and is also involved in competition dynamics between mobile
356 elements^{51,54}. Overall, our findings show that phylogroup B genomes are significantly larger and have a

357 wider pool of accessory genes than those from the other two phylogroups, possibly driven by the lower
358 prevalence of CRISPR-Cas systems in phylogroup B.

359

360 **AMR and defence systems are overrepresented in RGPs from phylogroups A and B**

361 We next assessed variation in the relative frequency of proteins encoded in RGPs from different
362 phylogroups. We observed that most functional categories are conserved across phylogroups. However,
363 proteins coding for replication, recombination and repair functions are more prevalent in phylogroups A
364 and B RGPs than those from phylogroup C (**Figure 5A**). Since these proteins are frequently involved in
365 mobilization, this finding may suggest that genomes in these phylogroups have more functional mobile
366 elements, with the ability to be horizontally transferred, while the RGPs in phylogroup C may be derived
367 from remnants of mobile elements that can no longer be mobilized.



368

369 **Figure 5.** Distribution of functional categories across RGPs and masked genomes from the different
 370 phylogroups. Bar and boxplots in blue represent phylogroup A, yellow B, and grey C. **A.** Relative
 371 frequencies of cluster of orthologous groups categories. The relative frequencies were calculated by
 372 dividing the absolute counts for each category by the total number of clustered proteins found in each of
 373 the six groups. Error bars indicate the degree of variation across each COG category from each
 374 phylogroup across RGPs and masked genomes. The functional categories are indicated by capital letters,
 375 including: A, RNA processing and modification; B, chromatin structure and dynamics; C, energy

376 production and conversion; D, cell cycle control and mitosis; E, amino acid metabolism and transport; F,
377 nucleotide metabolism and transport; G, carbohydrate metabolism and transport; H, coenzyme
378 metabolism; I, lipid metabolism; J, translation; K, transcription; L, replication, recombination and repair;
379 M, cell wall/membrane/envelop biogenesis; N, cell motility; O, post-translational modification, protein
380 turnover, chaperone functions; P, inorganic ion transport and metabolism; Q, secondary structure; R,
381 general functional prediction only; S, function unknown; T, signal transduction; U, intracellular
382 trafficking and secretion; V, defence mechanisms; W, extracellular structures; Z, cytoskeleton. **B**
383 Boxplots of the variation in the number of AMR genes, defence systems, and virulence genes found in
384 RGPs and masked genomes across the three phylogroups. Absolute counts of genes and systems were
385 normalized to RGP and masked genome sequence lengths in each strain. Values above 0.05 were
386 considered as non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq$
387 0.001 , and **** $p \leq 0.0001$.

388 We next assessed to what extent RGPs and masked genomes vary in prevalence of genes for three types
389 of functions, which are often encoded on MGEs, including virulence, defence systems, and AMR. Since
390 the cumulative size of all RGPs is substantially smaller than that of masked genomes (**Table S2**), the
391 number of virulence genes, defence systems, and AMR genes were normalized to the sequence length of
392 the RGPs and masked genomes for each strain. We observed that the gene prevalence for these functions
393 is conserved across masked genomes from different phylogroups, while they are unevenly distributed in
394 RGPs. (**Figure 5B**).

395 Two important virulence factors were only present in some genomes from phylogroup C, and absent from
396 the other two phylogroups (**Table S9**). These genes (*exlA* and *exlB*) encode hemolysins, and when
397 genomes from phylogroup C carry these genes, the typical type-III secretion system (T3SS) machinery
398 found in most bacteria (encoding the toxins ExoS, ExoY, ExoT, and ExoU) is absent from these genomes,
399 supporting previous reports that these are mutually exclusive⁵⁵. In agreement with previous findings¹⁴,
400 we further found that two important genes encoding T3SS effector proteins (*exoS* and *exoU*) were
401 unevenly distributed across the phylogroups: the *exoS* gene was pervasive among genomes from
402 phylogroup A (99.5%, 1524/1531) and the majority of phylogroup C strains (28/43), while the *exoU* gene
403 was overrepresented in genomes from phylogroup B (408/435) and nearly absent in genomes from the
404 other two phylogroups (**Table S9**). Surprisingly, we also found 23 genomes with the atypical
405 *exoS*⁺/*exoU*⁺ genotype, all belonging to phylogroup A (**Table S9**). A high frequency of this genotype has
406 recently been reported in patients from the Brazilian Amazon and Peruvian hospitals^{56,57}. As expected⁵⁸,
407 some virulence genes were exclusively found on RGPs (i.e., absent from masked genomes): flagellar-
408 associated proteins *fleI/flag*, *flgL*, *fliC* and *fliD*, as well as *wzy*, which codes for an O-antigen chain length
409 regulator. All these virulence genes were found in RGPs from both phylogroups A and B.

410 In agreement with the important role of MGEs as vectors for AMR genes in *P. aeruginosa*^{9,59}, we found
411 that AMR genes were overrepresented in RGPs from phylogroups A and B (**Figures 5B and S9**). We
412 then calculated the relative proportion of different AMR classes across RGPs from the three phylogroups,
413 revealing that most AMR classes were overrepresented across RGPs from phylogroup B (**Figure S10**).
414 This result is consistent with our finding that RGPs play a significant role in the larger genome sizes from

415 this phylogroup (**Figure S4**). Point mutations linked to resistance to beta-lactams and quinolones were
416 observed for all phylogroups (**Table S11**).

417 A wide array of defence systems with a patchy distribution in closely related and distantly related strains
418 was recently characterized in *P. aeruginosa*, suggesting high rates of horizontal gene transfer ⁴¹.
419 According to this hypothesis, we would expect to observe an abundance of defence systems in RGPs,
420 when compared with masked genomes. Similar to our results for AMR genes, we found that defence
421 systems are indeed overrepresented in RGPs from phylogroups A and B (**Figure 5B**). Defence systems
422 such as the globally distributed restriction-modification and CRISPR-Cas systems were common in RGPs
423 from both phylogroups. Some rarer systems such as cyclic-oligonucleotide-based anti-phage signalling
424 systems (CBASS) ⁶⁰, Zorya, Gabija, Druantia ⁶¹, abortive infection ⁶², and bacteriophage exclusion
425 (BREX) ⁶³ were also observed in RGPs from phylogroups A and B (**Figure S11** and **Table S12**). In
426 contrast, dGTPases were absent from both phylogroups. Finally, we also observed that AMR and defence
427 systems are overrepresented in specific MLST profiles, including the high-risk clones ST111 and ST233
428 (**Figure S12**) ¹. Our results revealed that AMR and defence systems are pervasive in RGPs from
429 phylogroups A and B, and the majority of AMR classes are overrepresented in RGPs from phylogroup B.

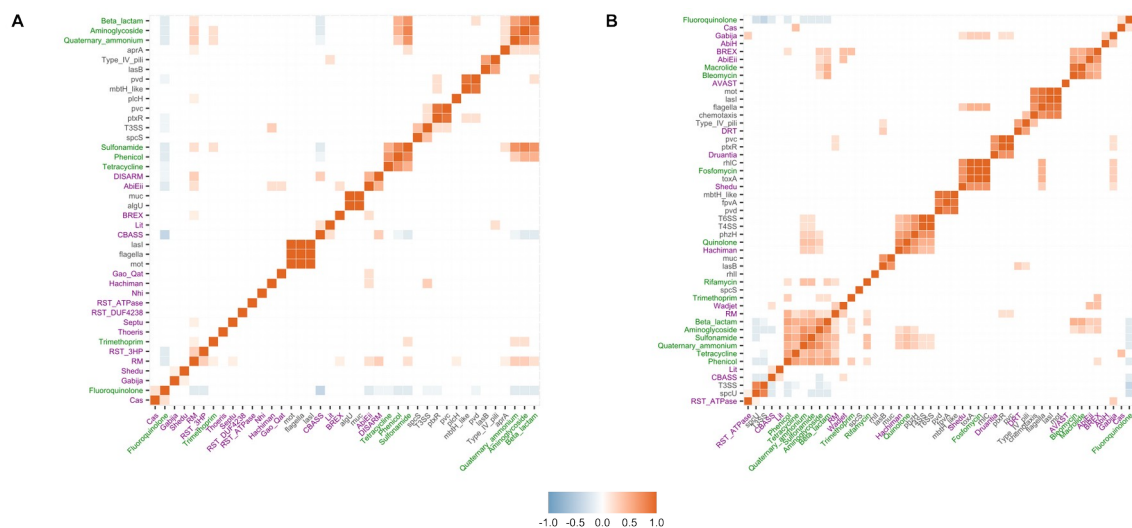
430

431 **AMR and defence systems are prevalent in ICEs/IMEs from phylogroups A and B**

432 Given that the distribution and clustering of defence systems in *P. aeruginosa* is not dependent on the
433 phylogenetic distance between all strains ⁴¹, and considering the high prevalence of ICEs/IMEs in this
434 species ⁶⁴, we explored the potential role of these elements as defence islands. To accurately detect these
435 MGEs, we focused our analysis on complete genomes. We noted that 12.6% of our collection consisted
436 of complete genomes (254/2009), including 172 genomes from phylogroup A, 78 from phylogroup B,
437 and 4 genomes from phylogroup C (**Table S2**). 215 out of the 254 complete genomes harboured a total of
438 477 ICEs and 76 IMEs (**Table S13**). These ICEs/IMEs were present in 136 genomes from phylogroup A,
439 77 from phylogroup B, and 2 from phylogroup C. Thus, ICEs/IMEs were pervasive in strains from
440 phylogroup B (77/78) and in the majority of strains from phylogroup A (136/172).

441 Nearly half of the ICEs/IMEs carried at least one AMR gene (228/553), with the ciprofloxacin-modifying
442 *crpP* gene and the sulphonamide-resistance *sull* gene being most frequent (**Table S14**). Indeed, the *crpP*
443 gene was recently shown to be widely dispersed across ICEs from *P. aeruginosa* ⁶⁵. Around one third of
444 the ICEs/IMEs (193/553) carried at least one defence system, resulting in a total of 250 defence systems
445 across the ICEs/IMEs and including 27 different types (**Figure S13** and **Table S14**). The most frequent
446 defence subtypes were CBASS-III and restriction-modification type-II ^{60,62}. Virulence genes were present
447 in a smaller proportion of the ICEs/IMEs (99/553) and showed higher variation in abundance across
448 ICEs/IMEs than AMR genes and defence systems do (**Figure S14**). The *exoU* gene encoding for the
449 effector protein and the *spcU* gene encoding for its chaperone were the most frequent virulence genes, all
450 in ICEs/IMEs from phylogroup B (**Table S14**).

451 We next explored to what extent the prevalence of these three functional groups is correlated across
452 ICEs/IMEs from the two larger phylogroups A and B. We observed that genes encoding resistance to
453 fluoroquinolones were negatively correlated with genes involved in resistance to other antibiotic classes,
454 and also with specific defence systems as restriction-modification and CBASS (**Figure 6A**). ICEs/IMEs
455 from phylogroup B carrying fluoroquinolone-encoding resistance genes were also negatively associated
456 with genes from the type-III secretion systems (**Figure 6B**). In contrast, genes encoding resistance to
457 distinct antibiotic classes (e.g., beta-lactams, aminoglycosides, and sulphonamides) were often positively
458 correlated in the ICEs/IMEs from both phylogroups, consistent with the previous observations that these
459 genes tend to be co-localized in genetic structures named integrons ⁶⁶. Virulence genes involved in
460 flagellar motility were also often correlated, either additionally with (phylogroup B) or without
461 (phylogroup A) genes involved in chemotaxis ⁶⁷. Defence systems BREX and AbiEii ^{62,63} were positively
462 correlated in ICEs/IMEs from phylogroup B. AMR and defence systems showed a high density in
463 ICEs/IMEs from phylogroups A and B, and their frequencies were positively correlated in both
464 phylogroups.



465

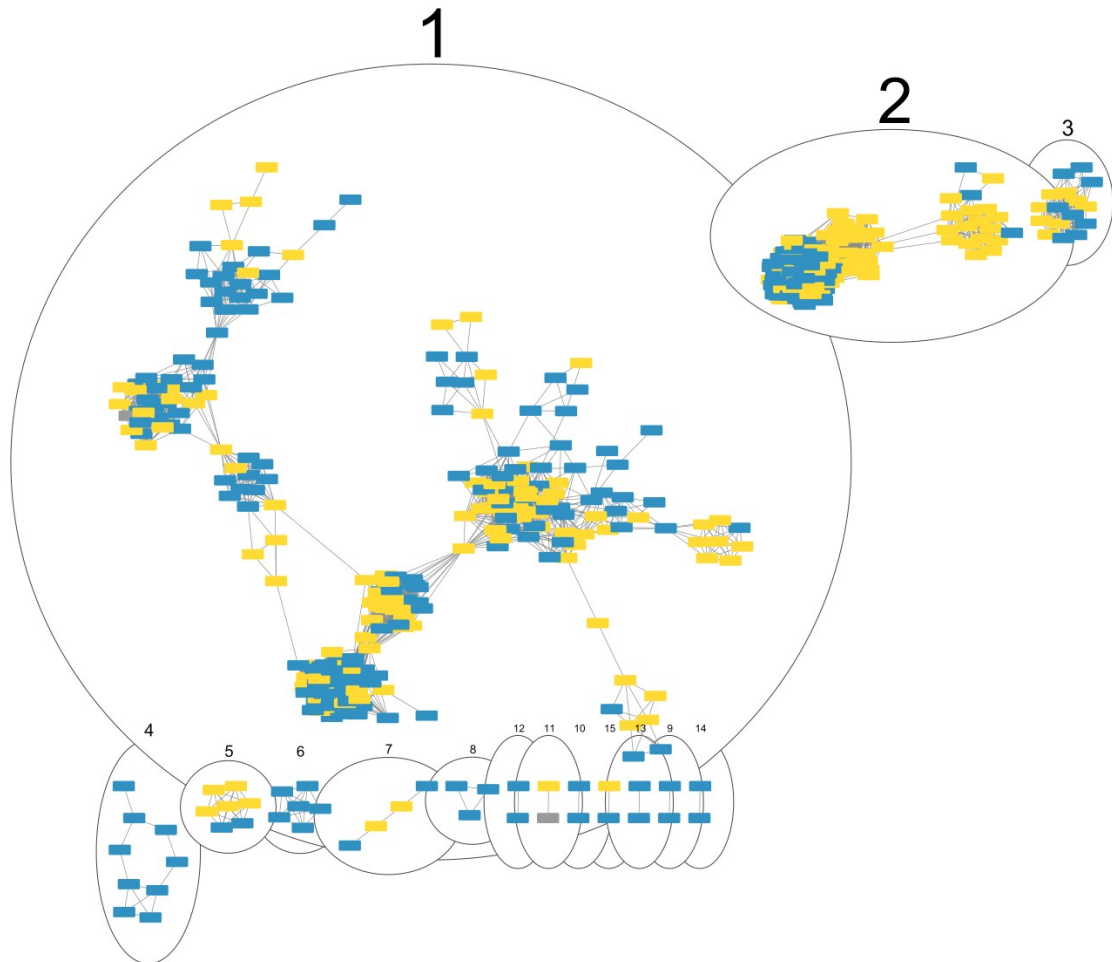
466 **Figure 6.** Correlation plots between AMR classes, virulence genes, and defence systems across
467 ICEs/IMEs from phylogroup A (**A**) and phylogroup B (**B**). The distribution of cargo genes across
468 ICEs/IMEs was converted into a presence/absence matrix. Correlation matrices were ordered using the
469 hierarchical clustering function. Positive correlations are shown in different shades of red, while negative
470 correlations are shown in different shades of blue. AMR genes and point mutations encoding resistance to
471 particular AMR classes are part of the AMRFinder database ³⁹, defence systems of defense-finder ⁴¹, and
472 virulence genes of the VFDB ⁴⁰. Virulence gene labels are coloured in black, AMR in green, and defence
473 systems in purple.

474

475 **ICEs/IMEs and RGPs from different phylogroups share high genetic similarity**

476 We next used an alignment-free sequence similarity comparison of the ICEs/IMEs to infer an undirected
477 network. The density plot showed a right-skewed distribution of pairwise distance similarities where the

478 vast majority of ICE/IME pairs shared little similarity, with a Jaccard Index value below 0.5 (**Figure**
479 **S15**), in accordance with the high diversity frequently observed across MGEs⁶⁸. To reduce the density
480 and increase the sparsity of the network, we used the mean Jaccard Index between all pairs of RGPs as a
481 threshold (0.12184). The network assigned 95.8% (530/553) of the ICEs/IMEs into 15 clusters (**Figure**
482 **7**). Nearly half of the ICEs/IMEs were grouped in cluster 1 (259/530, **Table S15**), which includes
483 representatives of the three phylogroups.



484

485 **Figure 7.** Network of clustered ICEs/IMEs from the three phylogroups, using the mean Jaccard Index
486 between all pairs of ICEs/IMEs as a threshold. Each ICE/IME is represented by a node, connected by
487 edges according to the pairwise distances between all ICE/IME pairs. Numbered ellipses represent
488 ICEs/IMEs that belong to the same cluster. The network has a clustering coefficient of 0.794, a density of
489 0.099, a centralization of 0.217, and a heterogeneity of 0.785. ICEs/IMEs from phylogroup A are
490 coloured in blue, from phylogroup B in yellow, and from phylogroup C in grey.

491 We then focused our analysis on the RGPs we extracted from all phylogroups (57901 RGPs in total). We
492 filtered out RGPs smaller than 10kb, and calculated the Jaccard Index between all pairwise of the
493 resulting 32744 RGPs. To reduce the density and increase the sparsity of the network, we used as a
494 threshold the mean value (0.0919429) of the estimated pairwise distances between the 32744 RGPs
495 identified in this study. The network assigned 99.7% (32651/32744) of the RGPs larger than 10kb into 51

496 clusters (**Figure S16**). While the majority of the RGP clusters were homogeneous for a given phylogroup,
497 we also observed DNA sharing events between different phylogroups. These findings suggest that RGPs
498 and ICEs/IMEs from different *P. aeruginosa* phylogroups share high genetic identity.

499

500 **Discussion**

501 In this work, we explored the pangenome of the opportunistic human pathogen *P. aeruginosa* in
502 consideration of its three main phylogroups. This approach allowed us to characterize defining properties
503 of each phylogroup. In particular, we identified genes that are prevalent in the small phylogroup C and
504 absent from members of the two larger phylogroups. These genes would have been classified to be part of
505 the accessory genome in conventional analyses of the pangenome of the species as a whole. In contrast,
506 our refined approach suggests that these genes have an evolutionary advantage in a specific genetic
507 context that is particular to this phylogroup⁶⁹. Moreover, phylogroup C is also clearly distinct from the
508 other two phylogroups A and B in having a significantly smaller genome size and a low relative
509 frequency of AMR and defence systems across RGPs. In addition, our results indicate an inverse
510 association between size of the phylogroup B accessory genome and presence of CRISPR-Cas systems.
511 This association could (but need not) be causal, such that a low prevalence of CRISPR-Cas defence
512 systems may possibly favour an increase in the size of the accessory genome. Remarkably, genomes
513 devoid of CRISPR-Cas systems in phylogroups A and B were generally significantly larger than those
514 with these systems, a trend that was no longer observed when only considering the non-RGP (“masked”)
515 genomes. This observation is consistent with the hypothesis that CRISPR-Cas systems can constrain
516 horizontal gene transfer in *P. aeruginosa*^{50–52,70}, at least for genomes belonging to the larger phylogroups.

517 The three phylogroups vary substantially in the distribution of AMR genes, defence systems, and
518 virulence genes. This variation is particularly apparent in the separate analyses of RGPs and masked
519 genomes. While the length of RGPs is substantially smaller than that of masked genomes, the absolute
520 counts of most defence systems were higher in RGPs than in masked genomes across the three
521 phylogroups (**Figure S11**). Curiously, representatives of the recently described set of defence systems
522 that are part of Doron’s seminal study⁶¹, such as Zorya, Wadjet, and Hachiman systems, were exclusively
523 found in RGPs across the three phylogroups. In Doron’s study, the authors demonstrated that the Wadjet
524 system provided protection against plasmid transformation in *Bacillus subtilis*, while the Zorya and
525 Hachiman systems mediated defence against bacteriophages. These findings highlight the important role
526 of defence systems encoded in RGPs in protecting genomes against infection by foreign DNA and their
527 contribution to MGE-MGE conflict. Moreover, AMR and defence systems are rare in RGPs from
528 phylogroup C, which may suggest that these strains are more often subjected to infection by foreign
529 DNA. Assuming that there is no sampling bias across the three phylogroups, then the smaller number of
530 phylogroup C members in public databases could thus be a consequence of the weaker arsenal of AMR
531 and defence systems. Alternatively, phylogroup C strains may indeed be underrepresented, for example if
532 they mainly occur in non-clinical habitats, which are usually less well sampled. Collecting *P. aeruginosa*

533 samples from distinct geographic regions and environments may further help us reconstruct variation in
534 metabolic competences and their connection to origin ⁷¹.

535 In general, our results underscore the role of ICEs/IMEs as vectors not only of AMR genes ⁵⁹, but also of
536 defence systems. Indeed, most of these systems show nonrandom clustering in defence islands and are
537 often co-localized with mobilome genes ^{61,72-74}. Co-occurrence of genes alone, however, does not infer an
538 ecological interaction between them ⁷⁵. Recently, it was proposed that the accessory genome of the genus
539 *Pseudomonas* is influenced by natural selection, showing a higher level of genetic structure than would be
540 expected if neutral processes governed the pangenome formation ⁷⁶. This suggests that coincident genes in
541 ICEs/IMEs are more likely to act together for the benefit of the host or to ensure their own maintenance
542 ^{9,11}. ICEs/IMEs, in particular, provide abundant material for the experimental study of bacterial defence
543 systems. For example, SXT ICEs in *Vibrio cholerae*, which are also involved in AMR, consistently
544 encode defence systems localized to a single hotspot of genetic shuffling ⁷⁷. Additionally, ICEs in
545 *Acidithiobacillia* carry type-IV CRISPR-Cas systems with remarkable evolutionary plasticity, which are
546 often involved in MGE-MGE warfare ⁷⁸. Moreover, a recent study proposed that size constraints may
547 account for the low abundance of large defence systems on prophages ⁴¹. In turn, the comparatively larger
548 size of ICEs/IMEs (when compared with prophages) ⁷⁹ may then explain that they commonly harbor large
549 systems such as BREX and defence island system associated with restriction–modification (DISARM) ⁸⁰
550 across our dataset (**Figure S13**). Even though the CBASS systems are not as prevalent as restriction-
551 modification and CRISPR-Cas systems across the bacterial phylogeny ⁴¹, three types of this system were
552 found across ICEs/IMEs from the larger phylogroups.

553 For our analyses, we used complete and draft genome assemblies retrieved from public databases.
554 However, incomplete genome assemblies likely impact RGP definition, due to highly fragmented
555 genomes, that might have inadvertently split RGPs into multiple contigs. With that in mind, we
556 subsampled the complete genomes from our collection and used these to accurately delineate ICEs/IMEs.
557 With the sequence similarity comparison between all pairs of ICEs/IMEs found in this study, as well as
558 between all pairs of RGPs, we were able to explore interactions between these elements, suggesting that
559 members of the same and of different phylogroups frequently undergo DNA shuffling events.
560 Importantly, this network-based approach using pairwise genetic distances of alignment-free *k*-mer
561 sequences between MGE pairs has bypassed the exclusion of non-coding elements, providing a more
562 comprehensive picture of MGE populations and dynamics ^{51,81}. Nevertheless, with the current progress in
563 sequence technology, especially including long-read sequencing, we envision a much larger number of
564 completely assembled *P. aeruginosa* genomes in the future, which will then improve reliable assessment
565 of the RGP composition and the role of particular MGEs or gene functions in shaping this species'
566 genome characteristics.

567 To conclude, our work used a refined approach to explore phylogroup-specific and pangenome dynamics
568 in *P. aeruginosa*. Members of phylogroup B contribute a comparatively larger number of pangenome
569 families, have larger genomes, and have a lower prevalence of CRISPR-Cas systems. AMR and defence
570 systems are pervasive in RGPs and ICEs/IMEs from phylogroups A and B, and these two functional
571 groups are often significantly correlated, including both positive and negative correlations. We also

572 observed multiple interaction events between the accessory genome content both between and within
573 phylogroups, suggesting that recombination events are frequent. Our conclusions are contingent on the
574 current range of sequenced genomes for *P. aeruginosa*. We cannot exclude that some groups, for example
575 phylogroup C and possibly its subgroups, are not fully represented in the currently available data. Future
576 sequencing efforts are likely to rectify such a possible problem, thus allowing to test the findings from our
577 study. Finally, our work provides a representative set of phylogenetically diverse *P. aeruginosa* strains,
578 the mPact strain panel, which should prove useful as a reference set for future functional analyses. Such
579 functional analyses may help to experimentally assess the underlying reasons for some of the correlations
580 identified in our study, for example the role of specific defence systems in RGP size expansion or in
581 mediating conflict between different MGE types.

582

583 **Declaration of interests**

584 We declare no competing interests.

585

586 **Contributors**

587 JB conceptualized and designed the work, acquired and analysed the data, interpreted the data, and wrote
588 the original draft of the manuscript. LT conceptualized and designed the work, acquired the data, and
589 interpreted the data. JF, FB, CU, JK, SN, and BT acquired and analysed the data. HS conceptualized and
590 designed the work, interpreted the data, and contributed to writing of the original draft of the manuscript.
591 HS, SN, and BT acquired funding for this work. All authors read, revised, and approved the final
592 manuscript.

593 **Data sharing**

594 Scripts for reproducing the analyses performed in this work are available at
595 https://gitlab.gwdg.de/botelho/pa_pangenome. The representative set of *P. aeruginosa* genomes and the
596 input file used for the network analysis in Figure 7 are available at the Figshare project
597 https://figshare.com/projects/P_aeruginosa_pangenome/155021. Analyses were made with a
598 combination of shell and R 4.0.3 scripting. Sequencing performed in this project were deposited in NCBI
599 under the Bioproject accession number PRJNA810040.

600

601 **Acknowledgments**

602 We acknowledge financial support from the German Science Foundation (grant SCHU 1415/12-2 to HS,
603 funding under Germany's Excellence Strategy EXC 2167-390884018 as well as the Research Training
604 Group 2501 TransEvo to HS and SN, and funding within the SFB 900 TP A2 to BT), the Max-Planck
605 Society (Max-Planck fellowship to HS), the Leibniz ScienceCampus Evolutionary Medicine of the Lung

606 (EvoLUNG, to HS and SN), and the BMBF program Medical Infection Genomics (AZ 0315827A to BT).
607 This work was supported by the Kiel Life Science Postdoc Award to JB and by the DFG Research
608 Infrastructure NGS_CC (project 407495230) as part of the Next Generation Sequencing Competence
609 Network (project 423957469). NGS was carried out at the Competence Centre for Genomic Analysis
610 (Kiel). This research was supported in part through high-performance computing resources available at
611 the Kiel University Computing Centre.

612

613 REFERENCES

- 614 1 Botelho J, Grosso F, Peixe L. Antibiotic resistance in *Pseudomonas aeruginosa* – Mechanisms,
615 epidemiology and evolution. *Drug Resist Updat* 2019; **44**: 100640.
- 616 2 De Oliveira DMP, Forde BM, Kidd TJ, *et al.* Antimicrobial resistance in ESKAPE pathogens.
617 *Clin Microbiol Rev* 2020; **33**. DOI:10.1128/CMR.00181-19.
- 618 3 Murray CJ, Ikuta KS, Sharara F, *et al.* Global burden of bacterial antimicrobial resistance in
619 2019: a systematic analysis. *Lancet* 2022; **399**: 629–55.
- 620 4 Horcajada JP, Montero M, Oliver A, *et al.* Epidemiology and treatment of multidrug-resistant and
621 extensively drug-resistant *Pseudomonas aeruginosa* infections. *Clin Microbiol Rev* 2019; **32**.
622 DOI:10.1128/CMR.00031-19/ASSET/3CDB73D6-9097-4D7E-A591-80EBAD992610/
623 ASSETS/GRAPHIC/CMR.00031-19-T004H.JPEG.
- 624 5 Tacconelli E, Carrara E, Savoldi A, *et al.* Discovery, research, and development of new
625 antibiotics: the WHO priority list of antibiotic-resistant bacteria and tuberculosis. *Lancet Infect*
626 *Dis* 2018; **18**: 318–27.
- 627 6 Koonin E V., Wolf YI. Genomics of bacteria and archaea: the emerging dynamic view of the
628 prokaryotic world. *Nucleic Acids Res* 2008; **36**: 6688–719.
- 629 7 Collins RE, Higgs PG. Testing the Infinitely Many Genes Model for the Evolution of the
630 Bacterial Core Genome and Pangenome. *Mol Biol Evol* 2012; **29**: 3413–25.
- 631 8 Arnold BJ, Huang I-T, Hanage WP. Horizontal gene transfer and adaptive evolution in bacteria.
632 *Nat Rev Microbiol* 2021; : 1–13.
- 633 9 Botelho J, Schulenburg H. The Role of Integrative and Conjugative Elements in Antibiotic
634 Resistance Evolution. *Trends Microbiol* 2020; **0**. DOI:10.1016/j.tim.2020.05.011.
- 635 10 Partridge SR, Kwong SM, Firth N, Jensen SO. Mobile Genetic Elements Associated with
636 Antimicrobial Resistance. *Clin Microbiol Rev* 2018; **31**: e00088-17.
- 637 11 Rocha Id EPC, Id DB. Microbial defenses against mobile genetic elements and viruses: Who
638 defends whom from what? *PLOS Biol* 2022; **20**: e3001514.
- 639 12 Pinilla-Redondo R, Russel J, Mayo-Muñoz D, *et al.* CRISPR-Cas systems are widespread
640 accessory elements across bacterial and archaeal plasmids. *Nucleic Acids Res* 2021; **1**: 13–4.
- 641 13 Horesh G, Taylor-Brown A, McGimpsey S, *et al.* Different evolutionary trends form the twilight
642 zone of the bacterial pan-genome. *Microb Genomics* 2021; **7**: 000670.
- 643 14 Ozer EA, Nnah E, Didelot X, Whitaker RJ, Hauser AR. The population structure of *Pseudomonas*
644 *aeruginosa* is characterized by genetic isolation of *exoU*⁺ and *exoS*⁺ lineages. *Genome Biol Evol*

- 645 2019; published online June 7. DOI:10.1093/gbe/evz119.
- 646 15 Trouillon J, Imbert L, Villard A-M, Vernet T, Attrée I, Elsen S. Determination of the two-
647 component systems regulatory network reveals core and accessory regulations across
648 *Pseudomonas aeruginosa* lineages. *Nucleic Acids Res* 2021; **1**: 13–4.
- 649 16 Trouillon J, Han K, Attrée I, Lory S, Kook H. The core and accessory Hfq interactomes across
650 *Pseudomonas aeruginosa* lineages. *Nat Commun* 2022 131 2022; **13**: 1–16.
- 651 17 Hilker R, Munder A, Klockgether J, *et al.* Interclonal gradient of virulence in the *Pseudomonas*
652 *aeruginosa* pangenome from disease and environment. *Environ Microbiol* 2015; **17**: 29–46.
- 653 18 Wiehlmann L, Cramer N, Tümmler B. Habitat-associated skew of clone abundance in the
654 *Pseudomonas aeruginosa* population. *Environ Microbiol Rep* 2015; **7**: 955–60.
- 655 19 Wiehlmann L, Wagner G, Cramer N, *et al.* Population structure of *Pseudomonas aeruginosa*.
656 *Proc Natl Acad Sci U S A* 2007; **104**: 8101–6.
- 657 20 Fischer S, Dethlefsen S, Klockgether J, Tümmler B. Phenotypic and Genomic Comparison of the
658 Two Most Common ExoU-Positive *Pseudomonas aeruginosa* Clones, PA14 and ST235.
659 *mSystems* 2020; **5**. DOI:10.1128/MSYSTEMS.01007-20.
- 660 21 Andrews S. FastQC A Quality Control tool for High Throughput Sequence Data.
661 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>. citeulike-article-id:11583827.
- 662 22 Babraham Bioinformatics - Trim Galore!
663 https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/ (accessed Feb 21, 2022).
- 664 23 Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies from
665 short and long sequencing reads. *PLOS Comput Biol* 2017; **13**: e1005595.
- 666 24 Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of *de novo* genome
667 assemblies: Fig. 1. *Bioinformatics* 2015; **31**: 3350–2.
- 668 25 Perrin A, Rocha EPC. PanACoTA: a modular tool for massive microbial comparative genomics.
669 *NAR Genomics Bioinforma* 2021; **3**. DOI:10.1093/nargab/lqaa106.
- 670 26 Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI
671 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* 2018; **9**:
672 5114.
- 673 27 Gautreau G, Bazin A, Gachet M, *et al.* PPanGGOLiN: Depicting microbial diversity via a
674 partitioned pangenome graph. *PLOS Comput Biol* 2020; **16**: e1007732.
- 675 28 Bazin A, Gautreau G, Médigue C, Vallenet D, Calteau A. panRGP: a pangenome-based method
676 to predict genomic islands and explore their diversity. *Bioinformatics* 2020; **36**: i651–8.
- 677 29 Minh BQ, Schmidt HA, Chernomor O, *et al.* IQ-TREE 2: New Models and Efficient Methods for
678 Phylogenetic Inference in the Genomic Era. *Mol Biol Evol* 2020; **37**: 1530–4.
- 679 30 Didelot X, Wilson DJ. ClonalFrameML: Efficient Inference of Recombination in Whole Bacterial
680 Genomes. *PLOS Comput Biol* 2015; **11**: e1004041.
- 681 31 Letunic I, Bork P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display
682 and annotation. *Nucleic Acids Res* 2021; **49**: W293–6.
- 683 32 Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic features.
684 *Bioinformatics* 2010; **26**: 841–2.

- 685 33 Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014; **30**: 2068–9.
- 686 34 Liu M, Li X, Xie Y, *et al.* ICEberg 2.0: an updated database of bacterial integrative and
687 conjugative elements. *Nucleic Acids Res* 2018; published online Nov 8.
688 DOI:10.1093/nar/gky1123.
- 689 35 Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the analysis
690 of massive data sets. *Nat Biotechnol* 2017 3511 2017; **35**: 1026–8.
- 691 36 Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. eggNOG-mapper v2:
692 Functional Annotation, Orthology Assignments, and Domain Prediction at the Metagenomic
693 Scale. *Mol Biol Evol* 2021; published online Oct 1. DOI:10.1093/MOLBEV/MSAB293.
- 694 37 Tatusov RL, Galperin MY, Natale DA, Koonin E V. The COG database: a tool for genome-scale
695 analysis of protein functions and evolution. *Nucleic Acids Res* 2000; **28**: 33–6.
- 696 38 Russel J, Pinilla-Redondo R, Mayo-Muñoz D, Shah SA, Sørensen SJ. CRISPRCasTyper:
697 Automated Identification, Annotation, and Classification of CRISPR-Cas Loci. *Cris J* 2020; **3**:
698 462–9.
- 699 39 Feldgarden M, Brover V, Haft DH, *et al.* Validating the AMRFinder Tool and Resistance Gene
700 Database by Using Antimicrobial Resistance Genotype-Phenotype Correlations in a Collection of
701 Isolates. *Antimicrob Agents Chemother* 2019; **63**. DOI:10.1128/AAC.00483-19.
- 702 40 Liu B, Zheng D, Jin Q, Chen L, Yang J. VFDB 2019: a comparative pathogenomic platform with
703 an interactive web interface. *Nucleic Acids Res* 2019; **47**: D687–92.
- 704 41 Tesson F, Hervé A, Mordret E, *et al.* Systematic and quantitative view of the antiviral arsenal of
705 prokaryotes. *Nat Commun* 2022 131 2022; **13**: 1–10.
- 706 42 Zhao X. BinDash, software for fast genome distance estimation on a typical personal laptop.
707 *Bioinformatics* 2019; **35**: 671–3.
- 708 43 Freschi L, Vincent AT, Jeukens J, *et al.* The *Pseudomonas aeruginosa* pan-genome provides new
709 insights on its population structure, horizontal gene transfer and pathogenicity. *Genome Biol Evol*
710 2018; published online Nov 29. DOI:10.1093/gbe/evy259.
- 711 44 Stover CK, Pham XQ, Erwin AL, *et al.* Complete genome sequence of *Pseudomonas aeruginosa*
712 PAO1, an opportunistic pathogen. *Nature* 2000; **406**: 959–64.
- 713 45 Roy PH, Tetu SG, Larouche A, *et al.* Complete genome sequence of the multiresistant taxonomic
714 outlier *Pseudomonas aeruginosa* PA7. *PLoS One* 2010; **5**.
715 DOI:10.1371/JOURNAL.PONE.0008842.
- 716 46 Morimoto Y, Tohya M, Aibibula Z, Baba T, Daida H, Kirikae T. Re-identification of strains
717 deposited as *Pseudomonas aeruginosa*, *Pseudomonas fluorescens* and *Pseudomonas putida* in
718 GenBank based on whole genome sequences. *Int J Syst Evol Microbiol* 2020; published online
719 Sept 16. DOI:10.1099/ijsem.0.004468.
- 720 47 Brockhurst MA, Harrison E, Hall JPJ, Richards T, McNally A, MacLean C. The Ecology and
721 Evolution of Pangenomes. *Curr Biol* 2019; **29**: R1094–103.
- 722 48 Filloux A. Protein secretion systems in *Pseudomonas aeruginosa*: An essay on diversity,
723 evolution, and function. *Front Microbiol* 2011; **2**: 155.
- 724 49 Koonin E V., Makarova KS, Wolf YI, Krupovic M. Evolutionary entanglement of mobile genetic

- 725 elements and host defence systems: guns for hire. *Nat Rev Genet* 2019; : 1–13.
- 726 50 Wheatley RM, MacLean RC. CRISPR-Cas systems restrict horizontal gene transfer in
727 *Pseudomonas aeruginosa*. *ISME J* 2020; published online Dec 21. DOI:10.1038/s41396-020-
728 00860-3.
- 729 51 Botelho J, Cazares A, Schulenburg H. The ESKAPE mobilome contributes to the spread of
730 antimicrobial resistance and CRISPR-mediated conflict between mobile genetic elements. *bioRxiv*
731 2022; : 2022.01.03.474784.
- 732 52 Pursey E, Dimitriu T, Paganelli FL, Westra ER, Houte S van. CRISPR-Cas is associated with
733 fewer antibiotic resistance genes in bacterial pathogens. *Philos Trans R Soc B* 2022; **377**.
734 DOI:10.1098/RSTB.2020.0464.
- 735 53 Pinilla-Redondo R, Mayo-Muñoz D, Russel J, *et al*. Type IV CRISPR-Cas systems are highly
736 diverse and involved in competition between plasmids. *Nucleic Acids Res* 2020; **48**: 2000–12.
- 737 54 León LM, Park AE, Borges AL, Zhang JY, Bondy-Denomy J. Mobile element warfare via
738 CRISPR and anti-CRISPR in *Pseudomonas aeruginosa*. *Nucleic Acids Res* 2021; **49**: 2114–25.
- 739 55 Reboud E, Basso P, Maillard AP, Huber P, Attrée I. Exolysin Shapes the Virulence of
740 *Pseudomonas aeruginosa* Clonal Outliers. *Toxins (Basel)* 2017; **9**.
741 DOI:10.3390/TOXINS9110364.
- 742 56 Horna G, Amaro C, Palacios A, Guerra H, Ruiz J. High frequency of the *exoU*⁺/*exoS*⁺ genotype
743 associated with multidrug-resistant ‘high-risk clones’ of *Pseudomonas aeruginosa* clinical isolates
744 from Peruvian hospitals. *Sci Rep* 2019; **9**. DOI:10.1038/S41598-019-47303-4.
- 745 57 Rodrigues YC, Furlaneto IP, Pinto Maciel AH, *et al*. High prevalence of atypical virulotype and
746 genetically diverse background among *Pseudomonas aeruginosa* isolates from a referral hospital
747 in the Brazilian Amazon. *PLoS One* 2020; **15**. DOI:10.1371/JOURNAL.PONE.0238741.
- 748 58 Arora SK, Bangera M, Lory S, Ramphal R. A genomic island in *Pseudomonas aeruginosa* carries
749 the determinants of flagellin glycosylation. *Proc Natl Acad Sci U S A* 2001; **98**: 9342–7.
- 750 59 Botelho J, Mourão J, Roberts AP, Peixe L. Comprehensive genome data analysis establishes a
751 triple whammy of carbapenemases, ICEs and multiple clinically relevant bacteria. *Microb*
752 *Genomics* 2020; : mgen000424.
- 753 60 Cohen D, Melamed S, Millman A, *et al*. Cyclic GMP–AMP signalling protects bacteria against
754 viral infection. *Nature* 2019; : 1–6.
- 755 61 Doron S, Melamed S, Ofir G, *et al*. Systematic discovery of antiphage defense systems in the
756 microbial pangenome. *Science (80-)* 2018; **359**: eaar4120.
- 757 62 Labrie SJ, Samson JE, Moineau S. Bacteriophage resistance mechanisms. *Nat Rev Microbiol*
758 2010 **85** 2010; **8**: 317–27.
- 759 63 Goldfarb T, Sberro H, Weinstock E, *et al*. BREX is a novel phage resistance system widespread
760 in microbial genomes. *EMBO J* 2015; **34**: 169–83.
- 761 64 Guglielmini J, Quintais L, Garcillán-Barcia MP, de la Cruz F, Rocha EPC. The Repertoire of ICE
762 in Prokaryotes Underscores the Unity, Diversity, and Ubiquity of Conjugation. *PLoS Genet* 2011;
763 **7**: e1002222.
- 764 65 Botelho J, Grosso F, Peixe L. ICEs Are the Main Reservoirs of the Ciprofloxacin-Modifying *crpP*

- 765 Gene in *Pseudomonas aeruginosa*. *Genes (Basel)* 2020; **11**: 889.
- 766 66 Ghaly TM, Geoghegan JL, Tetu SG, Gillings MR. The Peril and Promise of Integrons: Beyond
767 Antibiotic Resistance. *Trends Microbiol.* 2020; **28**: 455–64.
- 768 67 Matilla MA, Martín-Mora D, Gavira JA, Krell T. *Pseudomonas aeruginosa* as a Model To Study
769 Chemosensory Pathway Signaling. *Microbiol Mol Biol Rev* 2021; **85**.
770 DOI:10.1128/MMBR.00151-20.
- 771 68 Cury J, Oliveira PH, de la Cruz F, Rocha EPC. Host Range and Genetic Plasticity Explain the
772 Coexistence of Integrative and Extrachromosomal Mobile Genetic Elements. *Mol Biol Evol* 2018;
773 **35**: 2230–9.
- 774 69 Lassalle F, Muller D, Nesme X. Ecological speciation in bacteria: reverse ecology approaches
775 reveal the adaptive part of bacterial cladogenesis. *Res Microbiol* 2015; **166**: 729–41.
- 776 70 van Belkum A, Soriaga LB, LaFave MC, *et al.* Phylogenetic Distribution of CRISPR-Cas
777 Systems in Antibiotic-Resistant *Pseudomonas aeruginosa*. *MBio* 2015; **6**: e01796-15.
- 778 71 Saati-Santamaría Z, Baroncelli R, Rivas R, García-Fraile P. Comparative Genomics of the Genus
779 *Pseudomonas* Reveals Host- and Environment-Specific Evolution. *Microbiol Spectr* 2022;
780 published online Nov 10. DOI:10.1128/SPECTRUM.02370-22.
- 781 72 Makarova KS, Wolf YI, Snir S, Koonin E V. Defense Islands in Bacterial and Archaeal Genomes
782 and Prediction of Novel Defense Systems. *J Bacteriol* 2011; **193**: 6039–56.
- 783 73 Hussain FA, Dubert J, Elsherbini J, *et al.* Rapid evolutionary turnover of mobile genetic elements
784 drives bacterial resistance to phages. *Science (80-)* 2021; **374**: 488–92.
- 785 74 Vliet AHM van, Charity OJ, Reuter M. A *Campylobacter* integrative and conjugative element
786 with a CRISPR-Cas9 system targeting competing plasmids: a history of plasmid warfare? *Microb*
787 *Genomics* 2021; **7**: 000729.
- 788 75 Blanchet FG, Cazelles K, Gravel D. Co-occurrence is not evidence of ecological interactions.
789 *Ecol Lett* 2020; **23**: 1050–63.
- 790 76 Whelan FJ, Hall RJ, McInerney JO. Evidence for selection in the abundant accessory gene
791 content of a prokaryote pangenome. *Mol Biol Evol* 2021; published online May 7.
792 DOI:10.1093/molbev/msab139.
- 793 77 LeGault KN, Hays SG, Angermeyer A, *et al.* Temporal shifts in antibiotic resistance elements
794 govern phage-pathogen conflicts. *Science (80-)* 2021; **373**: eabg2166.
- 795 78 Moya-Beltrán A, Makarova KS, Acuña LG, *et al.* Evolution of Type IV CRISPR-Cas Systems:
796 Insights from CRISPR Loci in Integrative Conjugative Elements of Acidithiobacillia.
797 <https://home.liebertpub.com/crispr> 2021; published online Sept 28.
798 DOI:10.1089/CRISPR.2021.0051.
- 799 79 Cury J, Touchon M, Rocha EPC. Integrative and conjugative elements and their hosts:
800 composition, distribution and organization. *Nucleic Acids Res* 2017; **45**: 8943–56.
- 801 80 Ofir G, Melamed S, Sberro H, *et al.* DISARM is a widespread bacterial defence system with
802 broad anti-phage activities. *Nat Microbiol* 2017 31 2017; **3**: 90–8.
- 803 81 Acman M, van Dorp L, Santini JM, Balloux F. Large-scale network analysis captures biological
804 features of bacterial plasmids. *Nat Commun* 2020; **11**: 2452.

805

806

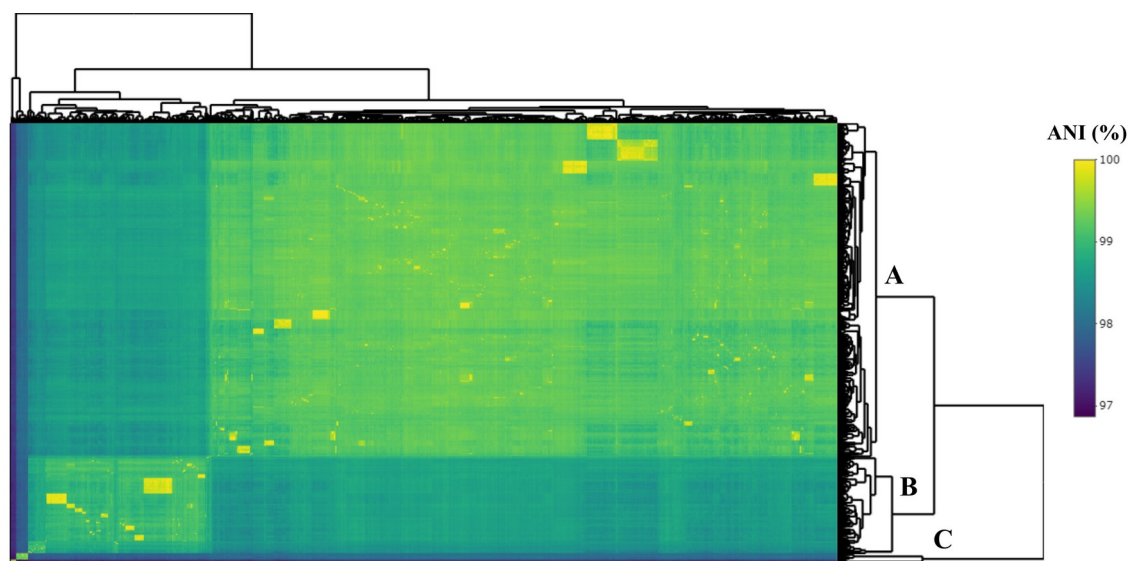
807

808 SUPPLEMENTARY FIGURES

809

810

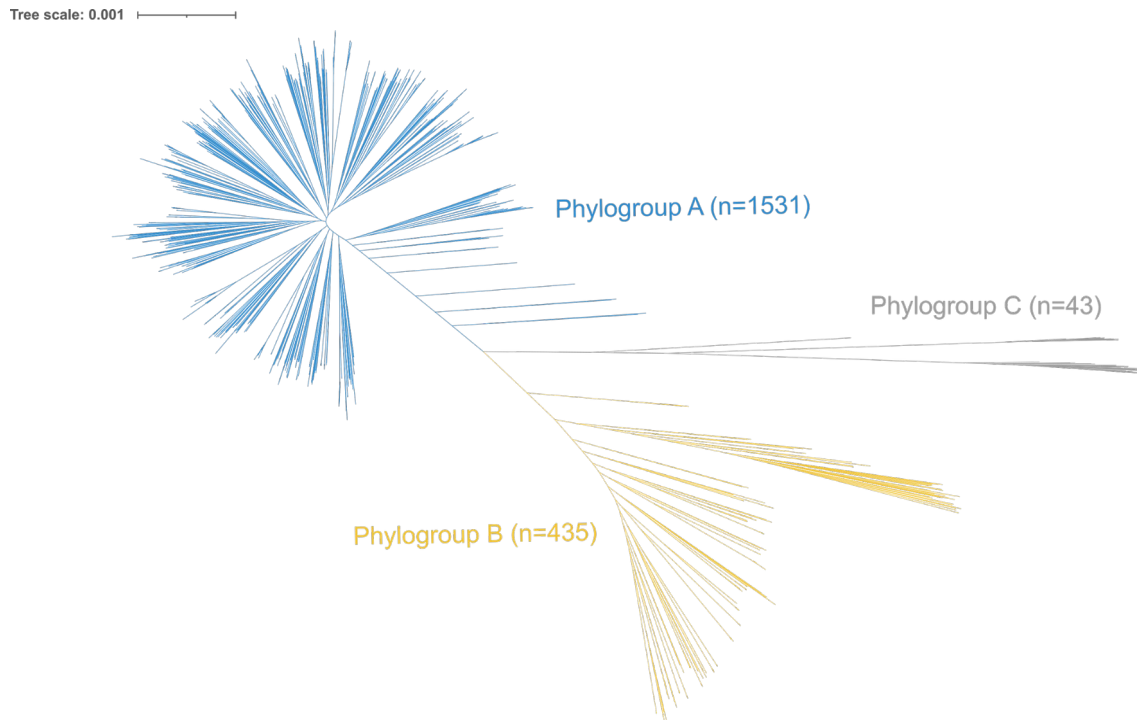
811



812

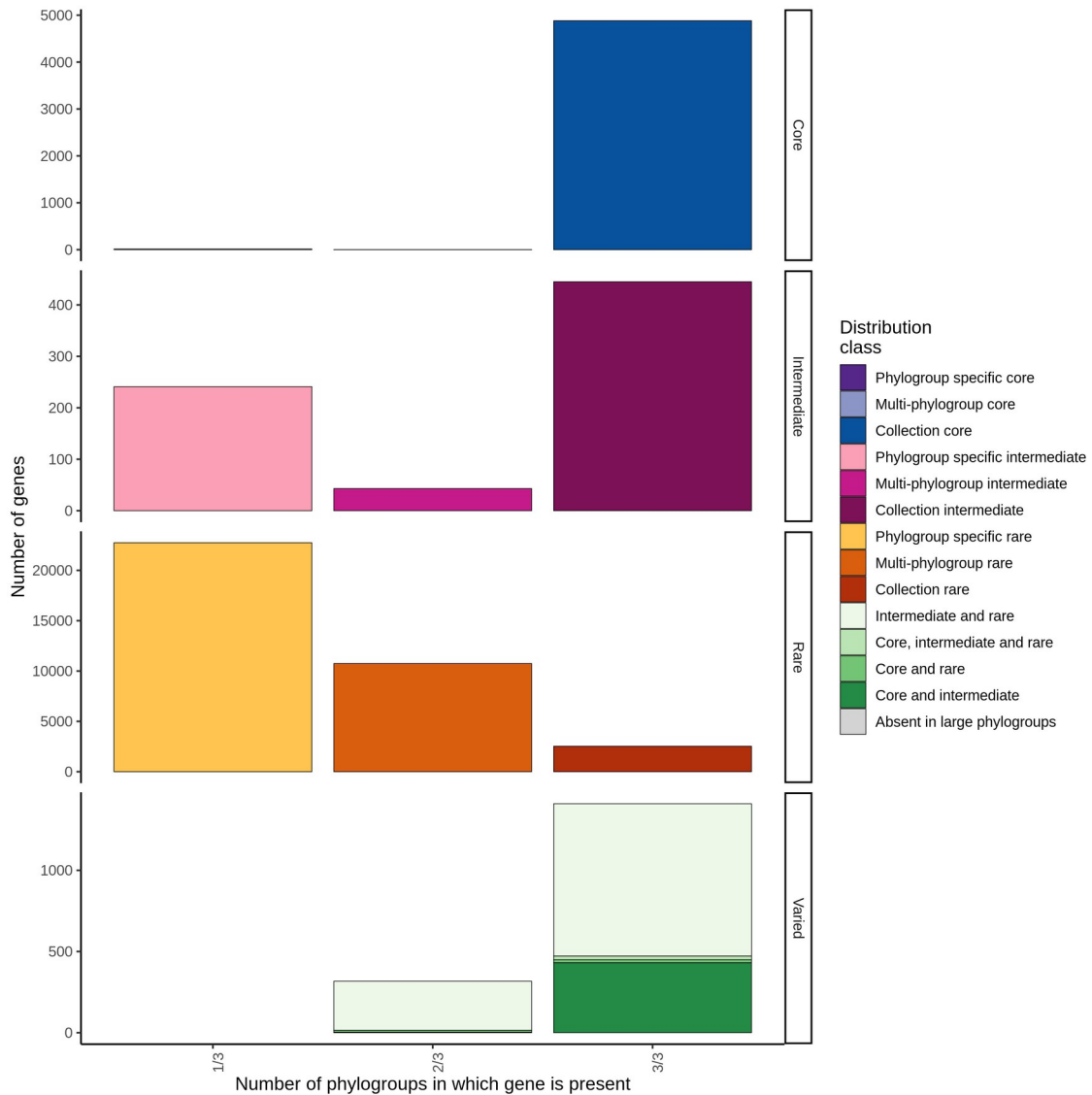
813 **Figure S1.** Matrix of Average Nucleotide Identity (ANI) scores between the 2009 *P. aeruginosa* genomes
814 used in this study. Row and column dendrograms are displayed. Hierarchical clustering was performed
815 with the complete-linkage clustering method. The three phylogroups identified in this study are
816 highlighted in the row dendrogram.

817



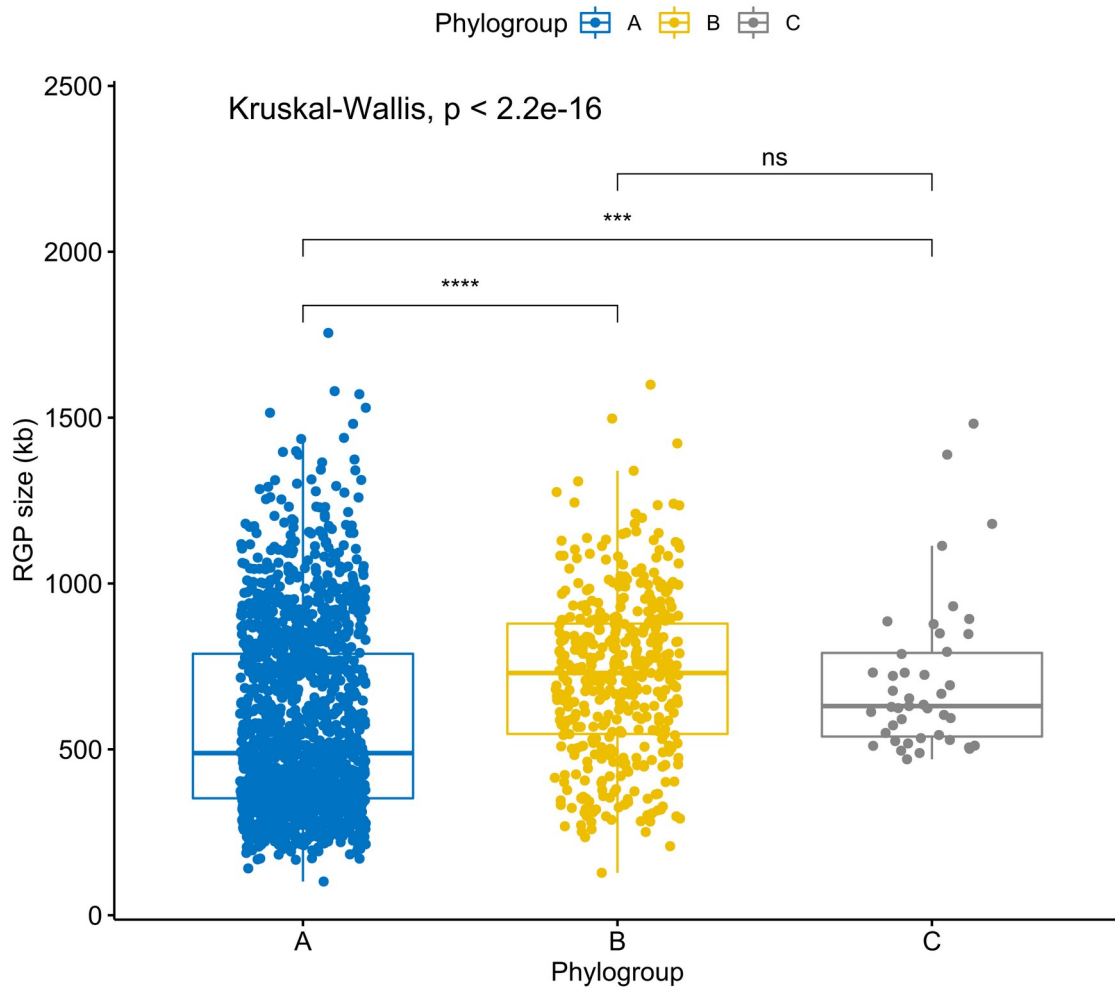
818

819 **Figure S2.** Maximum-likelihood tree of the softcore-genome alignment of all *P. aeruginosa* isolates used
820 in this study (n=2009), corrected for recombination. The scale bar represents the genetic distance.
821 Members of phylogroup A are coloured in blue, B in yellow, and C in grey.



822

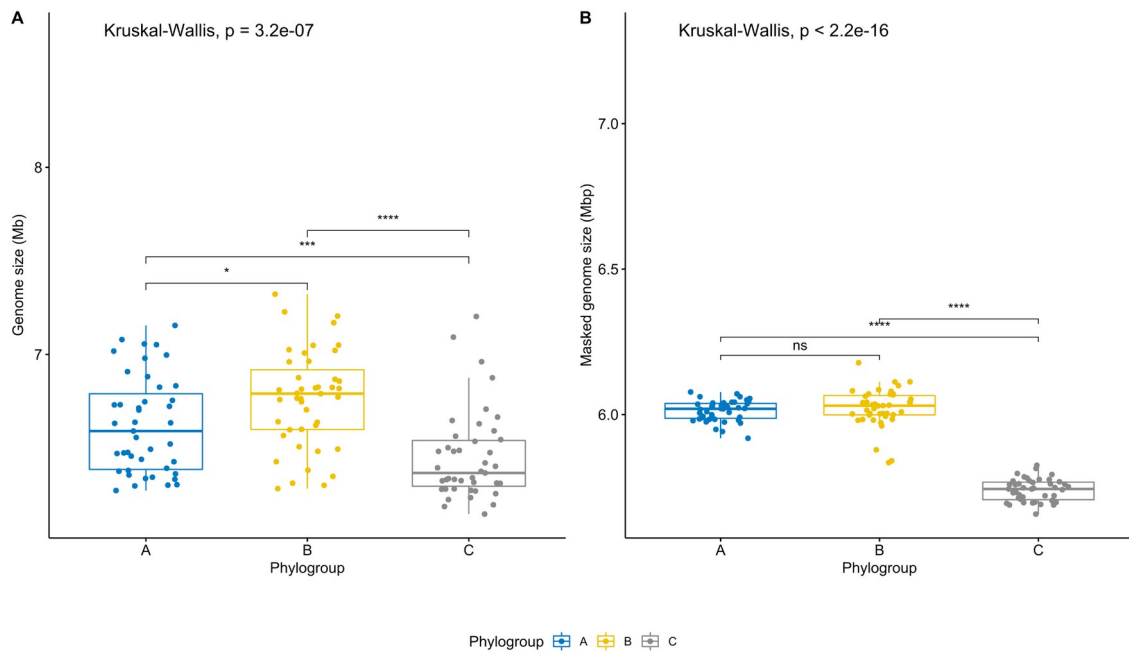
823 **Figure S3.** Barplots of the distribution of gene families into core, intermediate, rare, or varied parts of the
 824 pangenome across phylogroups. The first column shows genes that are specific to a given phylogroup,
 825 and further classified into core ($\geq 95\%$), intermediate, rare ($\leq 15\%$), or varied. The second column shows
 826 genes that are specific to two phylogroups, and their classification into core, intermediate, rare, or varied.
 827 The third column shows genes that are present across all three phylogroups, and their classification into
 828 core, intermediate, rare, or varied. A different colour is assigned to each classification. To create the plot,
 829 we modified the R script available in https://github.com/ghoresh11/twilight/blob/master/classify_genes.R.



830

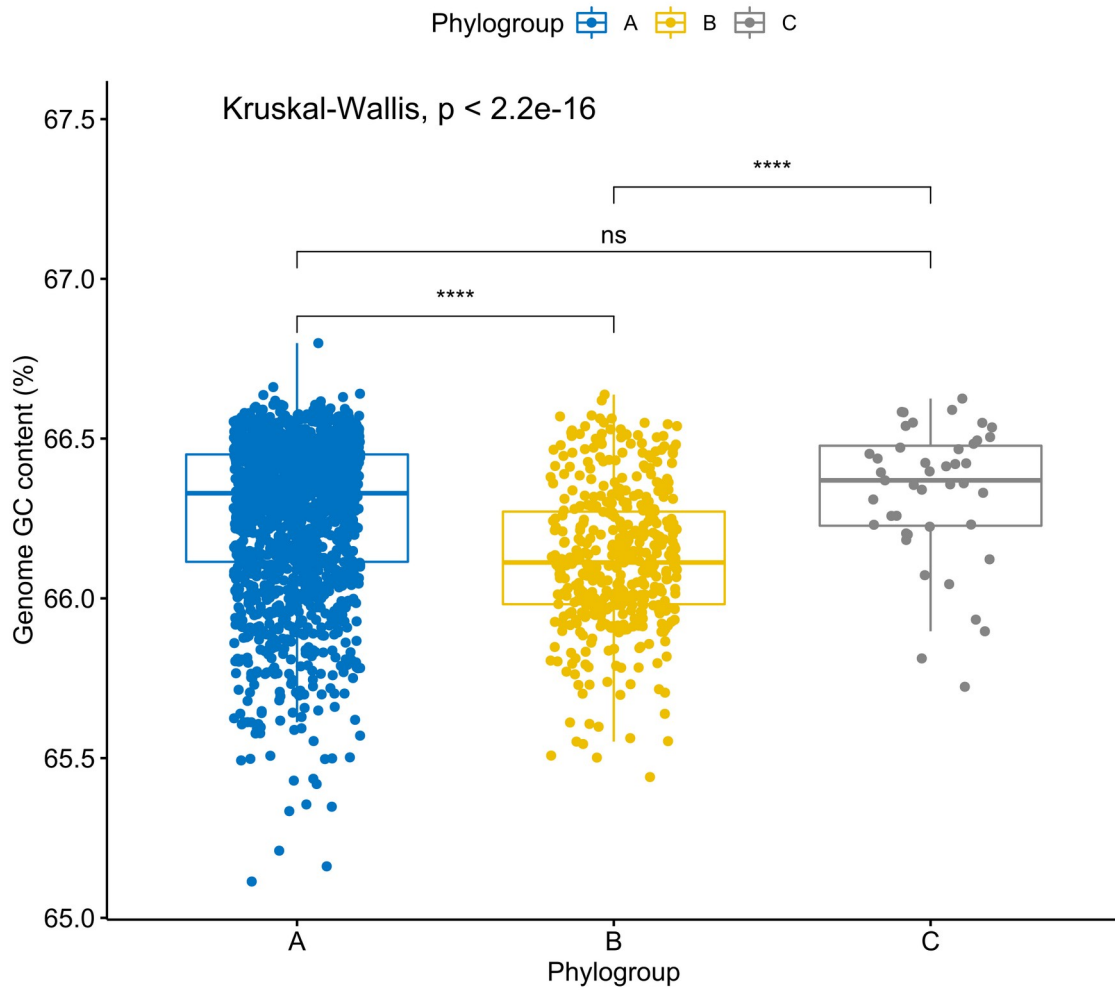
831

832 **Figure S4.** Boxplots showing the variation in RGP size across the three phylogroups. Values above 0.05
833 were considered as non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$, *** p
834 ≤ 0.001 , and **** $p \leq 0.0001$. Boxplots in blue represent phylogroup A, yellow B, and grey C.

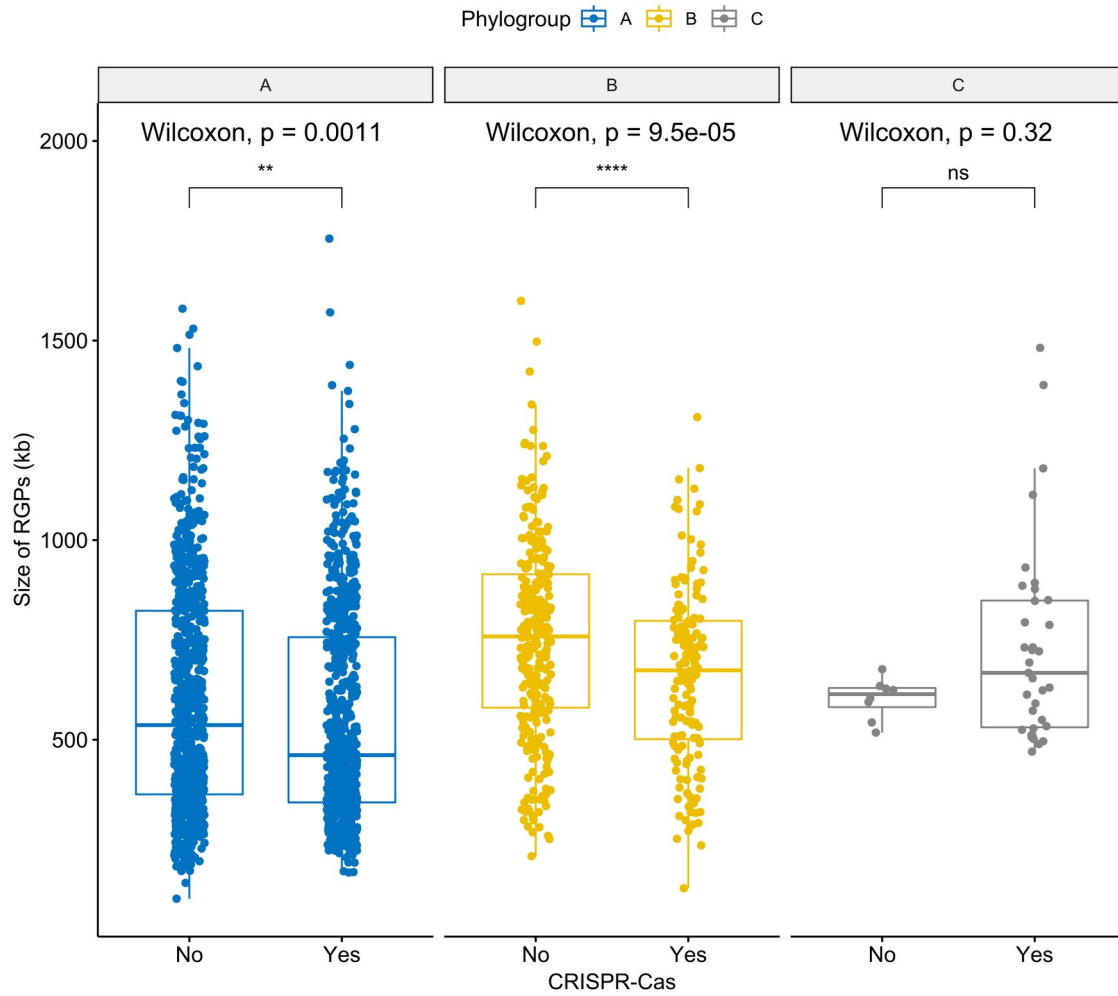


835

836 **Figure S5.** Boxplots representing the variation in genome size (A) and masked genome size (B) across
837 the three phylogroups. Phylogroup sample sizes were adjusted to sample size of the smallest group,
838 phylogroup C (with 43 genomes). Values above 0.05 were considered as non-significant (ns). Stars
839 indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq 0.001$, and **** $p \leq 0.0001$. Boxplots in
840 blue represent phylogroup A, yellow B, and grey C.



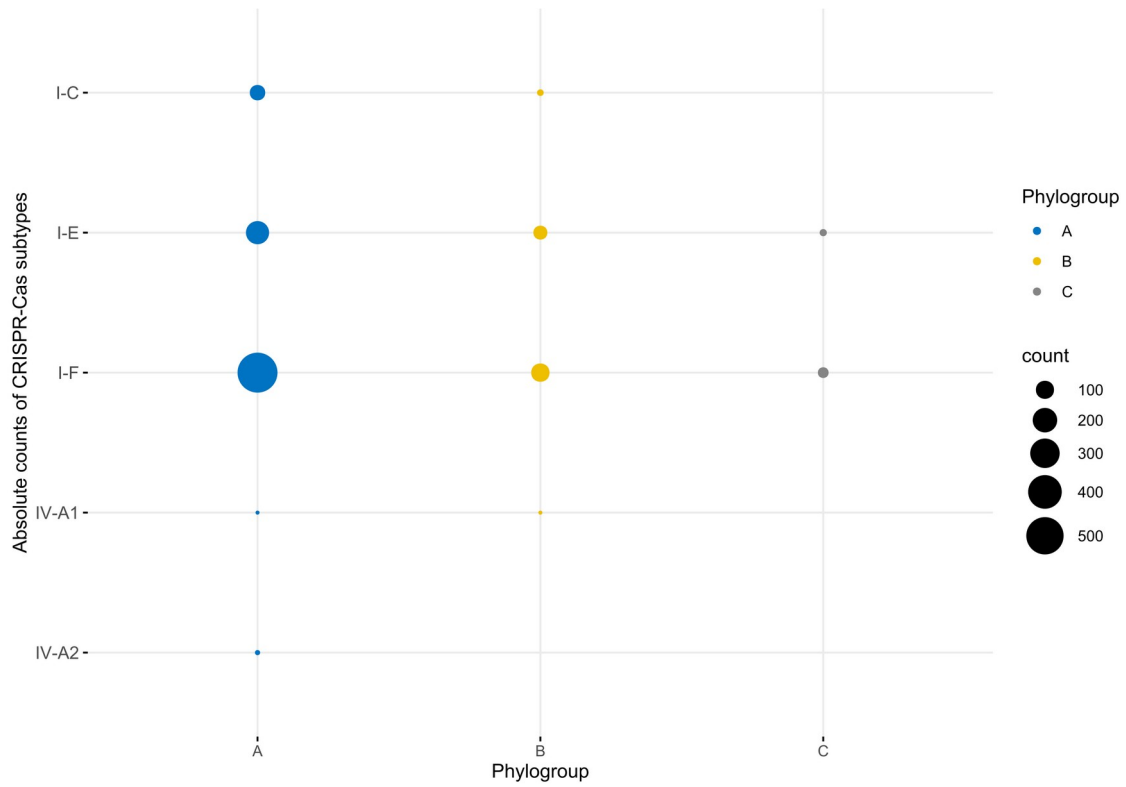
842 **Figure S6.** Boxplots showing the variation in GC content across the three phylogroups. Values above
843 0.05 were considered as non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$,
844 *** $p \leq 0.001$, and **** $p \leq 0.0001$. Boxplots in blue represent phylogroup A, yellow B, and grey C.



845

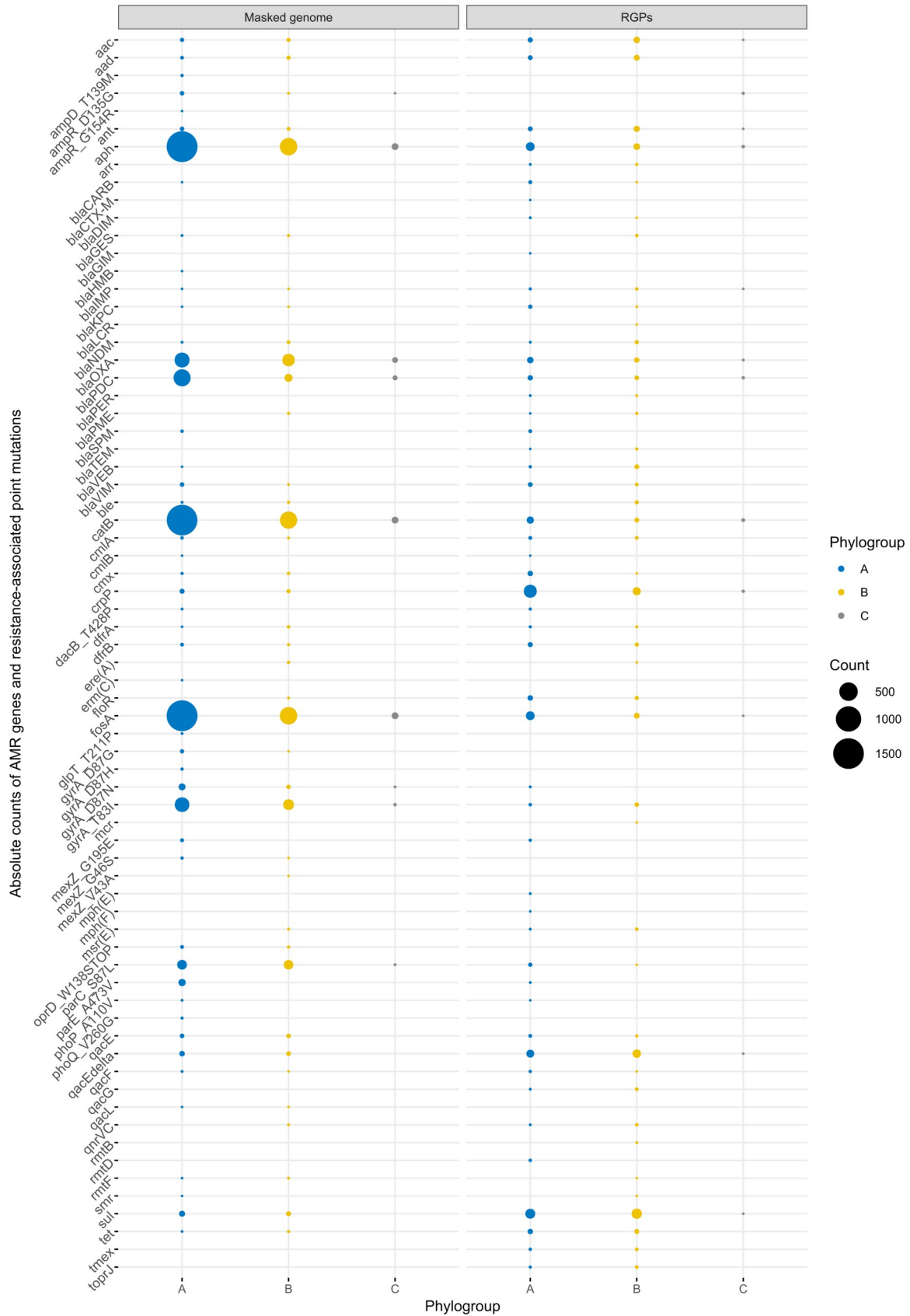
846 **Figure S7.** Boxplots representing the variation in the size of RGPs across pairs of conspecific genomes
847 from the same phylogroup with and without CRISPR-Cas systems. Values above 0.05 were considered as
848 non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq 0.001$, and
849 **** $p \leq 0.0001$. Boxplots in blue represent phylogroup A, yellow B, and grey C.

850



851

852 **Figure S8.** Absolute counts of CRISPR-Cas subtypes identified across genomes from the three
853 phylogroups. Circles in blue represent phylogroup A, yellow B, and grey C. Circle size is proportional to
854 the number of absolute counts.



855

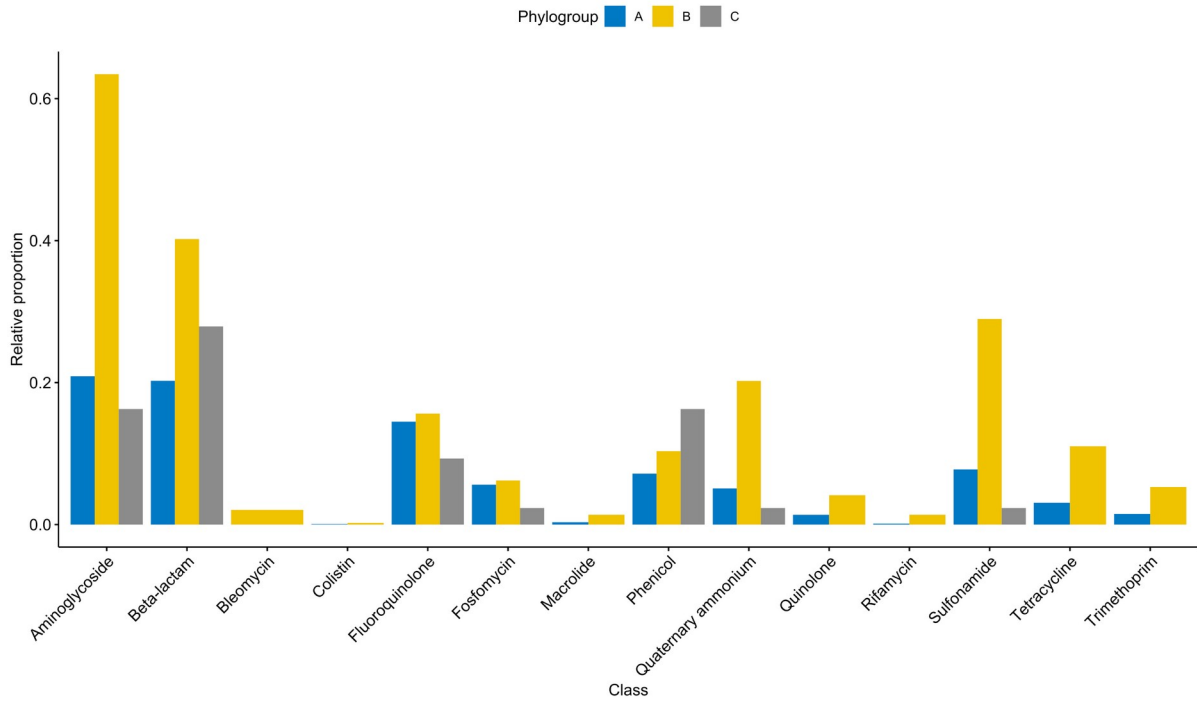
856 **Figure S9.** Absolute counts of AMR genes and resistance-associated point mutations across masked
 857 genomes and RGPs from the three phylogroups. Genes and mutations are part of the AMRFinder

1
 2
 3

36

858 database (45). Circle size is proportional to the number of absolute counts. Circles in blue represent
859 phylogroup A, yellow B, and grey C.

860

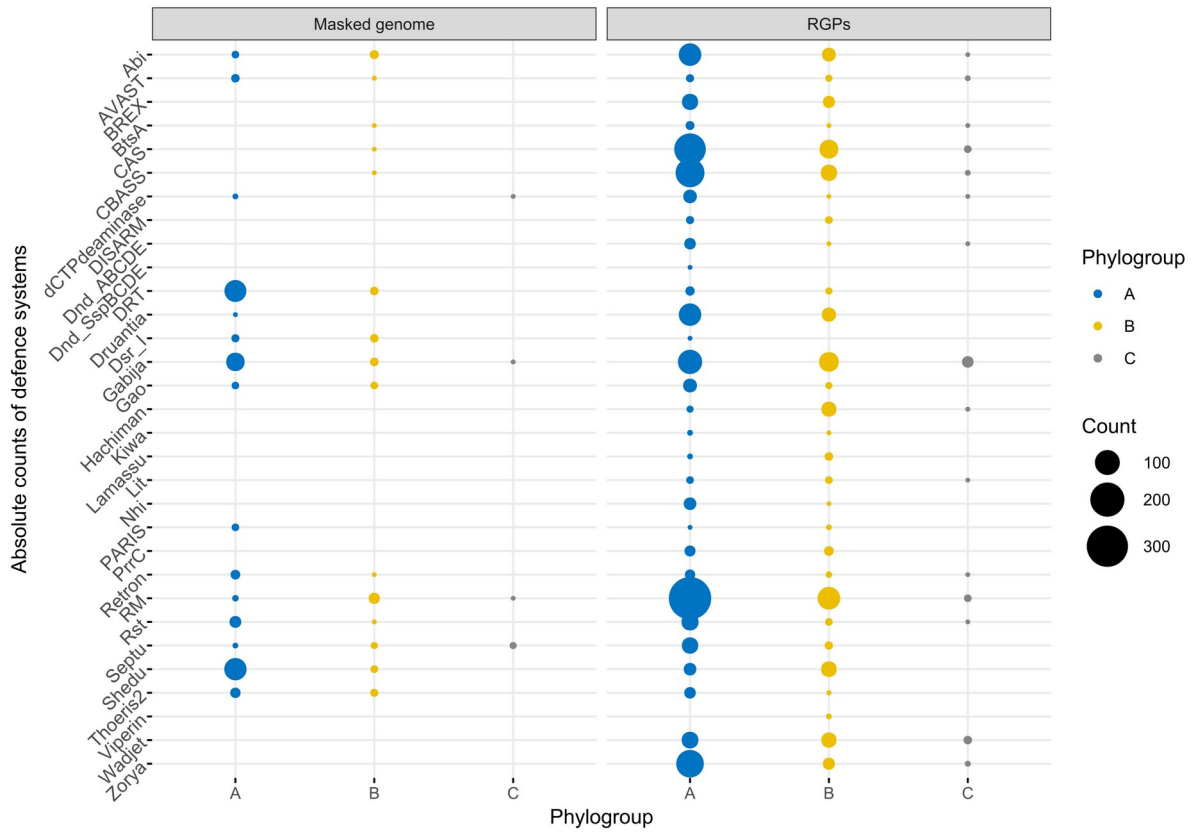


861

862 **Figure S10.** Barplots showing the relative proportion of genes encoding resistance to antibiotics from
863 different classes across RGPs from the three phylogroups. Genes were normalized to the total number of
864 genomes found in each phylogroup. Bars in blue represent phylogroup A, yellow B, and grey C.

865

866



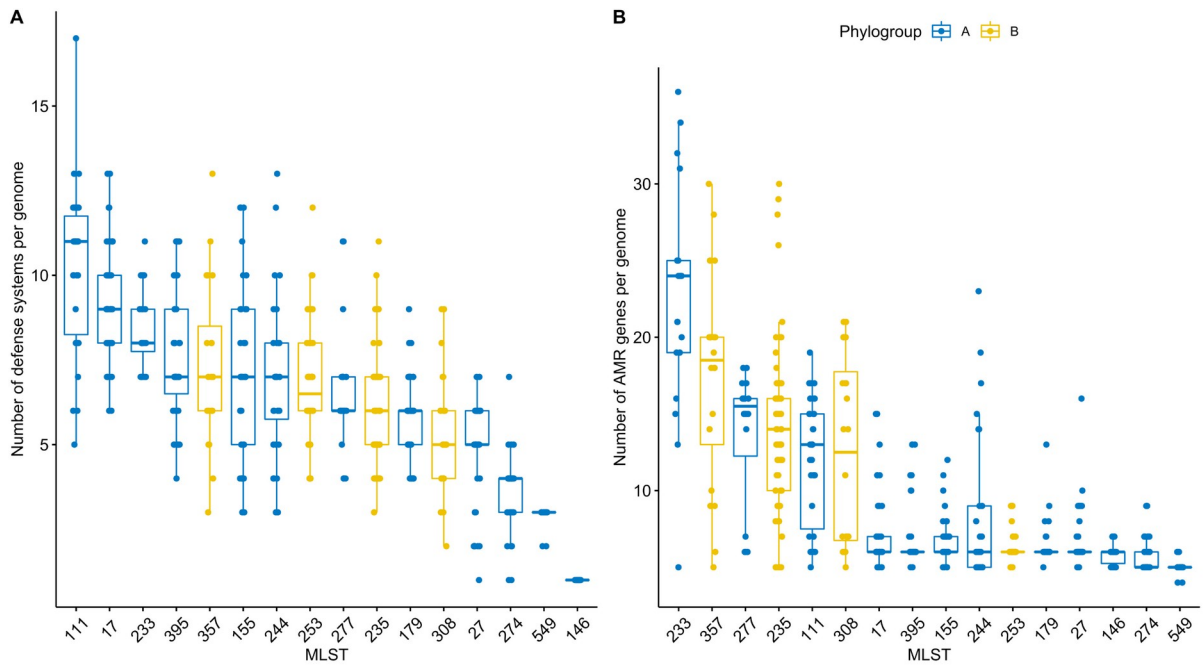
867

868 **Figure S11.** Absolute counts of defence systems across masked genomes and RGP from the three
 869 phylogroups. Defence systems are part of the defense-finder database (36). Circle size is proportional to
 870 the number of absolute counts. Circles in blue represent phylogroup A, yellow B, and grey C.

871

872

873



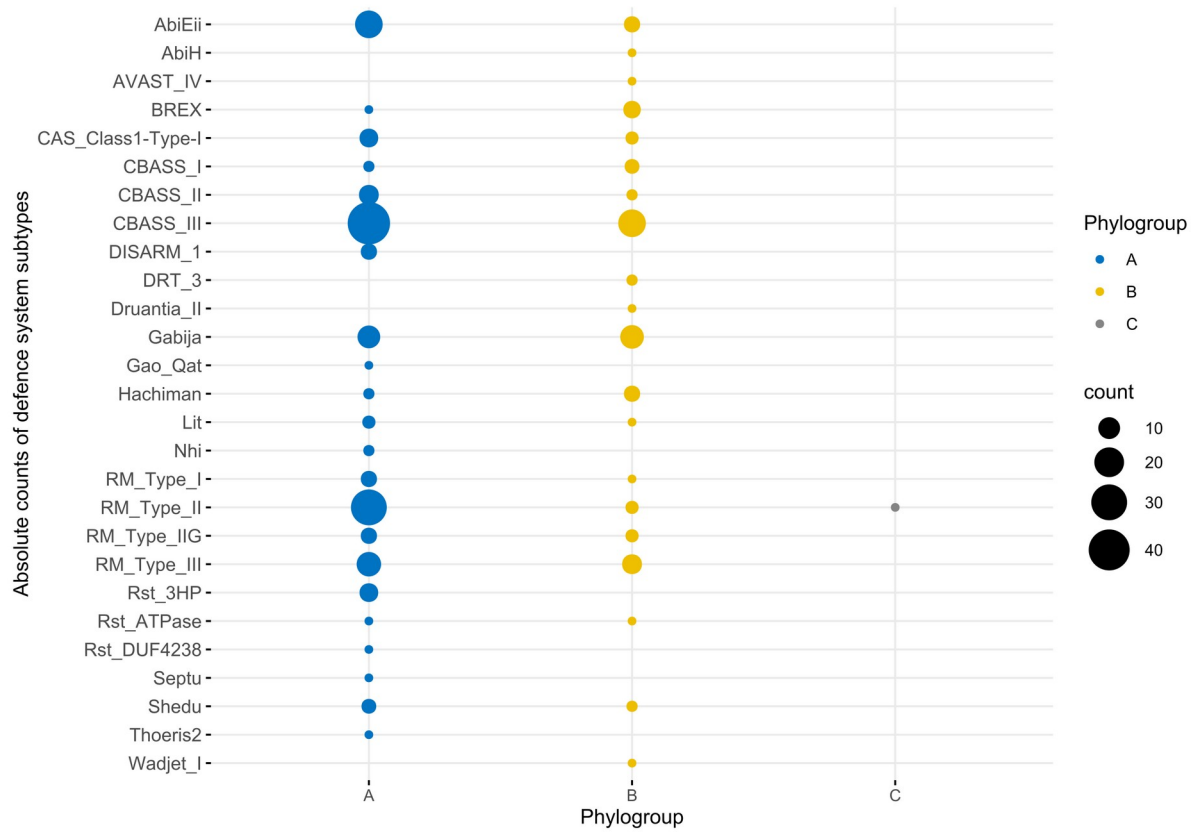
874

875 **Figure S12.** Boxplots representing the variation in the number of defense systems **(A)** and AMR genes
876 **(B)** found across the genomes from the main MLST profiles found in this study. Only MLST profiles
877 with at least 10 genomes are shown. Boxplots are ordered in descending order by the median values.

878

879

880

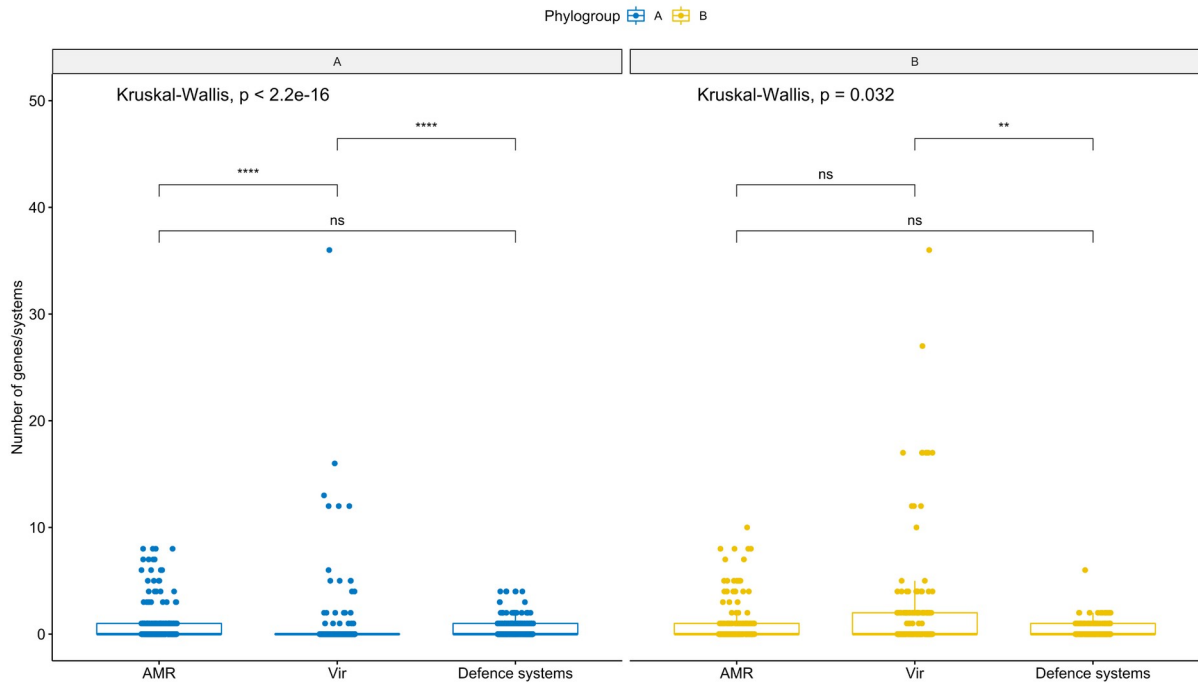


881

882 **Figure S13.** Absolute counts of defence systems across ICEs/IMEs from the three phylogroups. Defence
 883 systems are part of the defense-finder database (36). Circle size is proportional to the number of absolute
 884 counts. Circles in blue represent phylogroup A, yellow B, and grey C.

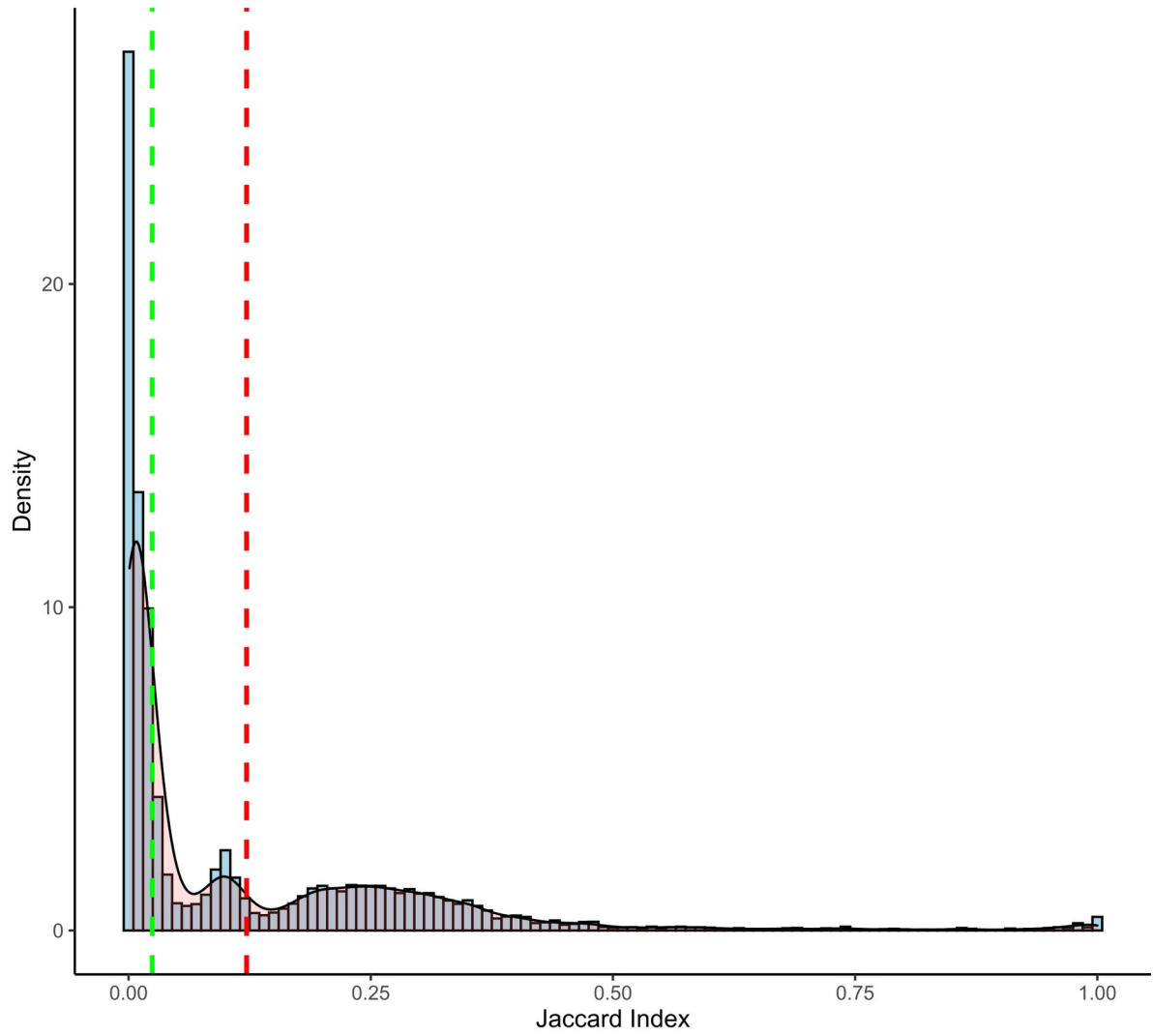
885

886



887

888 **Figure S14.** Boxplots representing the variation in the number of AMR genes, defence systems, and
889 virulence genes found in ICEs/IMEs across the two larger phylogroups A and B. Values above 0.05 were
890 considered as non-significant (ns). Stars indicate significance level: * $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq$
891 0.001, and **** $p \leq 0.0001$. Boxplots in blue represent phylogroup A, yellow B, and grey C.



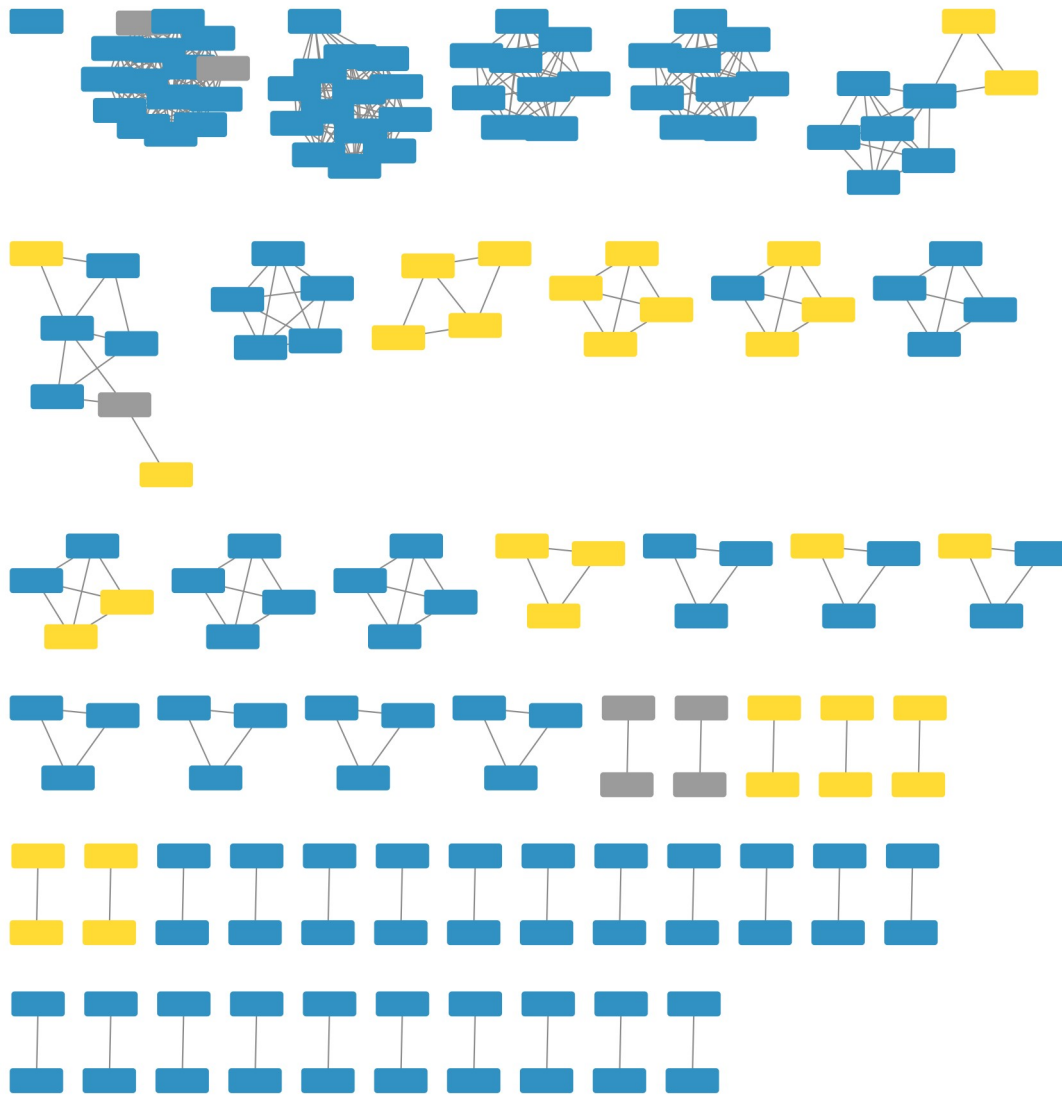
892

893

894 **Figure S15.** Histogram showing the right-skewed distribution of the Jaccard Index between all pairs of
895 ICEs/IMEs. Median and mean values are highlighted by vertical dashed lines in green and red,
896 respectively.

897

898



899

900 **Figure S16.** Network of clustered RGPs from the three phylogroups, using the mean Jaccard Index
901 between all pairs of RGPs as a threshold. Each RGP is represented by a node, connected by green edges
902 according to the pairwise distances between all RGPs pairs. Numbered ellipses represent RGPs that
903 belong to the same cluster. The network has a clustering coefficient of 0.777, a density of 0.007, a
904 centralization of 0.026, and a heterogeneity of 0.755. RGPs from phylogroup A are coloured in blue, from
905 phylogroup B in yellow, and from phylogroup C in grey.

906

907

908