1    **Genomic diversity of hospital-acquired infections revealed through prospective whole**

2    **genome sequencing-based surveillance**

3

4    Mustapha M. Mustapha[1,2]*, Vatsala R. Srinivasa[1,2]*, Marissa P. Griffith[1,2], Shu-Ting Cho[1], Daniel

5    R. Evans[1], Kady Waggle[1,2], Chinelo Ezeonwuka[1,2], Daniel J. Snyder[3], Jane W. Marsh[1,2], Lee H.

6    Harrison[1,2], Vaughn S. Cooper[3], Daria Van Tyne[1,^]

7

8    [1]Division of Infectious Diseases, University of Pittsburgh School of Medicine, Pittsburgh,

9    Pennsylvania, USA

10   [2] Microbial Genomic Epidemiology Laboratory, Center for Genomic Epidemiology, University of

11   Pittsburgh, Pittsburgh, Pennsylvania, USA

12   [3] Department of Microbiology and Molecular Genetics, and Center for Evolutionary Biology and

13   Medicine, University of Pittsburgh School of Medicine, Pennsylvania, USA

14   *These authors contributed equally

15   ^Correspondence to: vantyne@pitt.edu

16

17   **Keywords:** Whole genome sequencing, pangenome, antimicrobial resistance, horizontal gene

18   transfer, evolution

**Abstract**

Healthcare-associated infections (HAIs) cause mortality, morbidity, and waste of healthcare resources. HAIs are also an important driver of antimicrobial resistance, which is increasing around the world. Beginning in November 2016, we instituted an initiative to detect outbreaks of HAI using prospective whole genome sequencing-based surveillance of bacterial pathogens collected from hospitalized patients. Here we describe the biodiversity of bacteria sampled from hospitalized patients at a single center, as revealed through systematic analysis of their genomes. We sequenced the genomes of 3,004 bacterial isolates from hospitalized patients collected over a 25-month period. We identified bacteria belonging to 97 distinct species, which were distributed among 14 species groups. Within these groups, isolates could be distinguished from one another by both average nucleotide identity (ANI) and principal component analysis of accessory genes (PCA-A). Genetic distances between isolates and rates of evolution varied between different species, which has implications for the selection of distance cut-offs for outbreak analysis. Antimicrobial resistance genes and the sharing of mobile genetic elements between different species were frequently observed. Overall, this study describes the population structure of pathogens circulating in a single healthcare setting, and shows how investigating microbial population dynamics can inform genomic epidemiology studies.

**Importance**

Hospitalized patients are at increased risk of becoming infected with antibiotic-resistant organisms. We used whole-genome sequencing to survey and compare over 3,000 bacterial isolates collected from hospitalized patients at a large medical center over a two-year period. We identified nearly 100 different bacterial species, suggesting that patients can be infected with a wide variety of different organisms. When we examined how genetic relatedness differed between species, we found that different species are likely evolving at different rates within our hospital. This is significant because the identification of bacterial outbreaks in the hospital

45 currently relies on genetic similarity cut-offs, which are often applied uniformly across

46 organisms. Finally, we found that antibiotic resistance genes and mobile genetic elements were

47 abundant among the bacterial isolates we sampled. Overall, this study provides an in-depth

48 view of the genomic diversity and evolution of bacteria sampled from hospitalized patients, as

49 well as genetic similarity estimates that can inform hospital outbreak detection and prevention

50 efforts.

51

52 **Background**

53 Healthcare-associated infections (HAIs) affect over half a million people in the United States

54 each year, and annual direct hospital costs for treating HAIs are estimated at over $30 billion[1-3].

55 A relatively small number of bacterial species account for the majority of the burden of antibiotic-

56 resistant HAIs. Organisms belonging to the ESKAPE (*Enterococcus faecium*, *Staphylococcus*

57 *aureus*, *Klebsiella pneumoniae*, *Acinetobacter baumannii*, *Pseudomonas aeruginosa* and

58 *Enterobacter* spp.) group of pathogens are particularly problematic, due to their high burden of

59 HAIs and frequent multidrug resistance[2,4]. In addition, while *Clostridioides difficile* is not highly

60 antibiotic resistant, toxin-producing *C. difficile* lineages associated with significant patient

61 morbidity and mortality have emerged in recent years, making this organism an urgent health

62 threat[5].

63　　　Healthcare institutions such as hospitals and long-term care facilities constitute a unique

64 ecological niche for the proliferation and spread of antibiotic-resistant pathogens. The hospital

65 environment has a constant flow of vulnerable populations, and widespread use of antimicrobial

66 medications and cleaning agents provide selective pressure for the emergence and expansion

67 of drug-resistant bacterial strains[6]. Likewise, pathogens causing HAIs possess several common

68 biological traits that facilitate their survival and spread in healthcare environments. These traits

69 include frequent presence and acquisition of antimicrobial resistance, asymptomatic carriage,

70 and the ability to survive for prolonged periods on environmental surfaces such as medical

71 equipment, or in water systems[7-9]. These factors make healthcare settings a key contributor to

72 the increase of antibiotic-resistant bacterial infections worldwide.

73 Epidemiologic surveillance of HAIs requires timely and accurate ascertainment of strain

74 type to identify patients infected with genetically related strains of the same pathogen.

75 Surveillance using whole genome sequencing (WGS) is the gold standard for the detection of

76 outbreaks, and has provided significant insight into the population structure of hospital-

77 associated bacterial infections[10,11]. To improve the detection of hospital-associated transmission

78 at our medical center, we began conducting prospective WGS surveillance of clinical bacterial

79 isolates from hospitalized patients in November 2016, with the aim of identifying previously

80 undetected outbreaks and characterizing pathogen transmission routes. Our approach, called

81 Enhanced Detection of Hospital-Associated Transmission (EDS-HAT), combines prospective

82 bacterial WGS surveillance with data mining of the electronic health record to identify outbreaks,

83 including those that would otherwise go undetected, and their transmission routes[12-15]. In

84 conducting this work, we have collected and sequenced the genomes of thousands of bacterial

85 isolates. Systematic analysis of the genomes of these isolates can increase our understanding

86 of the diversity of bacteria causing HAIs[16].

87 Here we describe the genomic diversity, evolutionary rates, antimicrobial resistance

88 gene repertoires, and mobile genetic elements carried by over 3,000 bacterial isolates sampled

89 from patients at an academic medical center over 25 months. We uncovered a large and

90 diverse number of species causing HAIs at our center, and showed how different population

91 structures and evolutionary rates among these species can impact epidemiologic studies.

92 Systematic analyses of antimicrobial resistance genes and mobile genetic elements revealed

93 both species-specific differences as well as broader trends, and uncovered new avenues for

94 further investigation.

95

96 **Results**

**Pangenome analysis highlights the diversity of bacteria causing HAIs**

97

98  The objective of this study was to use WGS to examine the genetic diversity of HAIs at a single

99  medical center over a multi-year period, and to understand how this diversity impacts genomic

100 epidemiology and outbreak investigations. A total of 3,004 bacterial isolates collected from

101 2,046 unique patients at the University of Pittsburgh Medical Center (UPMC) from November

102 2016 through November 2018 were sequenced and analyzed. Isolates were distributed among

103 14 species groups (Supplementary Tables 1 and 2, Fig. 1). The largest proportion of isolates

104 were sampled from the respiratory tract (33.4%) followed by urinary tract (20.6%), tissue/wound

105 (20.6%), stool (16.7%, all *C. difficile*), and blood (8.7%) (Fig. 1). The distribution of isolated

106 species was similar between blood and tissue/wound, while the urinary tract, respiratory tract,

107 and stool samples had different species compositions. *P. aeruginosa* was the most prevalent

108 species isolated, with 863 isolates (28.7% of all isolates) collected from 653 unique patients.

109 Other prevalent species included toxin-producing *C. difficile* (16.7%), methicillin-resistant *S.*

110 *aureus* (MRSA, 14%) and vancomycin-resistant *E. faecalis* and *E. faecium* (VRE, 8.2%). The

111 remaining ten species groups contained less than 200 isolates each (Supplementary Table 1).

112 Genome sizes were highly variable, and ranged from a median length of 2.9Mb for MRSA to

113 7.6Mb for *Burkholderia* spp. (Fig. 2a). Pangenome collection curves constructed for genera

114 containing multiple species showed that *Citrobacter* spp. and *Acinetobacter* spp. had the

115 greatest pangenome diversity, perhaps due to the large number of different species sampled for

116 these groups (Fig. 2b, Supplementary Table 2). Pangenome collection curves for individual

117 species showed large differences in pangenome diversity between species (Fig. 2c), with MRSA

118 and VRE *faecium* genomes having the lowest diversity, while *P. aeruginosa, C. freundii,* and *S.*

119 *marcescens* had the greatest pangenome diversity of all species collected. The large and open

120 pangenome of *P. aeruginosa* is well known[17], however the pangenome diversity of *C. freundii*

121 and *S. marcescens* are not well described.

**Differences in bacterial population structures revealed by average nucleotide identity (ANI) and accessory gene content analysis**

Analysis of ANI and accessory genome contents are useful methods for assigning bacterial species, as well as understanding bacterial population structures[18-20]. Because the species of each isolate collected by the EDS-HAT project was initially assigned by the clinical microbiology laboratory, we first conducted pairwise comparisons of ANI for all isolate genomes, plus additional reference genomes downloaded from the NCBI database, and used a standard 95% ANI cut-off to group genomes into the same or different species[18]. This method resulted in the identification of 97 different species among the collected isolates (Supplementary Table 2). An example of ANI-based classification of *Citrobacter* spp. is shown in Fig. 3a. As expected, several species groups were highly diverse and were composed of multiple different species, including *Acinetobacter* spp., *Burkholderia* spp., *Citrobacter* spp., *Providencia* spp., *Pseudomonas* spp., and *Stenotrophomonas* spp. (Fig. 3a, Supplementary Fig. 1). Several other species groups, such as ESBL-producing *Klebsiella* spp., *Proteus* spp. and *Serratia* spp., were composed of one dominant species (*K. pneumoniae, P. mirabilis,* and *S. marcescens*), and a small number of isolates belonging to other species (Supplementary Table 1). ANI analysis of *P. aeruginosa* identified 15 isolates (1.7% of all *P. aeruginosa* collected) that belonged to a different species and could be clearly separated from the rest of the *P. aeruginosa* population by ANI (Supplementary Fig. 2). These 15 isolates all had greater than 95% ANI with the Group 3 PA7 genome[21], indicating that they belonged to this divergent group of *P. aeruginosa*. Overall, these findings highlight the potential discordance between species assignment based on clinical laboratory testing versus genome sequence analysis.

While ANI measures nucleotide identity in regions that are shared between two genomes, the accessory genes, which by definition are variably present in different genomes, can also be used to identify differences between bacterial species[42,43]. We constructed principal component analysis plots based on accessory gene content (PCA-A) for species groups

148    containing multiple species and with multiple isolates represented (Fig. 3b, Supplementary Fig.

149    1). The PCA-A plot for *Citrobacter* spp. isolates was largely congruent with species clustering by

150    ANI (Fig. 3b), and the same was true for *Acinetobacter* spp. and *Stenotrophomonas* spp. as well

151    (Supplementary Fig. 1). The *S. marcescens* isolates we collected could be clearly separated

152    into five different clades by both ANI and PCA-A; we arbitrarily named these clades A-E

153    (Supplementary Table 1, Supplementary Fig. 3). We observed that the pairwise ANI distribution

154    among all *S. marcescens* isolates included comparisons of isolates in different clades that fell

155    below the 95% ANI threshold used to distinguish species from one another (Fig. 3c,

156    Supplementary Fig. 2). Isolates within each *S. marcescens* clade always shared greater than

157    95% ANI with isolates in at least one other clade, however comparisons of isolates in Clade A

158    with isolates in either Clade C or Clade E fell below the 95% ANI threshold for same-species

159    comparisons (Supplementary Fig. 3). PCA-A clearly separated these clades from one another

160    (Fig. 3c), suggesting that each clade possessed a unique set of clade-specifying genes

161    (Supplementary Table 3). These data suggest that the *S. marcescens* population we sampled

162    may be in the process of diverging into distinct sub-species.

163    We also explored whether PCA-A could be used to cluster isolates belonging to different

164    genetic lineages within a single species (Fig. 3e-g). We analyzed isolates belonging to the

165    dominant lineages of toxin-producing *C. difficile* (Fig. 3e), VRE *faecium* (Fig. 3f), and MRSA

166    (Fig. 3g), and found in all cases that PCA-A could generally separate isolates belonging to

167    different STs. *C. difficile* isolates belonging to ST1, ST2, ST8, and ST42 were clearly separated

168    from one another (Fig. 3e). *E. faecium* isolates belonging to ST736 were clearly separated from

169    isolates belonging to ST17, ST18, and ST1471, which showed some overlap with one another

170    (Fig. 3f). Finally, MRSA isolates belonging to ST8 were clearly separated from isolates

171    belonging to ST5 and ST105, however the latter STs (which belong to the same clonal complex)

172    were not distinguishable from one another (Fig. 3g). Analysis of gene enrichment among these

173    different STs revealed ST-specific gene repertoires, which were largely composed of predicted

174    mobile element genes and hypothetical proteins (Supplementary Tables 4-6). These data

175    suggest that analysis of variable gene content may be a useful complement to SNP-based

176    methods in epidemiologic investigations.

177    **Genetic diversity and evolutionary rates vary by species**

178    The EDS-HAT project was designed to detect genetically and epidemiologically connected

179    isolates sampled from different patients, and has successfully identified dozens of clusters

180    containing isolates that share common exposures or transmission chains[14,15,22]. In addition, a

181    significant number of patients in this study were repeatedly sampled. To understand how

182    genetic diversity varied by species, we compared within-patient, within-cluster, and between-

183    patient diversity for six different species by calculating pairwise SNP distances for all isolate

184    pairs belonging to the same ST (Fig. 4a).  In all cases, SNP differences for pairs of isolates

185    collected from the same patient were on average lower than those for pairs of isolates collected

186    from different patients, suggesting that patients were persistently colonized or infected with the

187    same bacterial strain that was repeatedly sampled. Despite only comparing isolates belonging

188    to the same ST, some same-patient comparisons for *P. aeruginosa* resulted in hundreds or

189    thousands of SNPs, which could reflect reinfection with a different strain or the presence of

190    hypermutator strains. Within-cluster comparisons were comparable to within-patient

191    comparisons, demonstrating that clustered isolates were also highly genetically related to one

192    another. We also found that there were substantial differences in median SNP distances

193    between different species, with *C. difficile* isolates having the lowest median pairwise SNPs

194    among isolates from the same patient (2 SNPs), and *P. aeruginosa* having the highest (15

195    SNPs). These data likely reflect the different genome sizes, as well as the different biology of

196    the organisms studied here, and have broader implications for the selection of SNP cut-offs for

197    the purposes of epidemiologic investigation.

198            We next compared the evolutionary rates of the *C. difficile*, VRE, MRSA, and *P.*

199    *aeruginosa* populations that we sampled. We used TreeTime[23] to estimate the nucleotide

200    substitution rates for the most frequently observed STs for each species (Fig. 4b,

201    Supplementary Table 7). Consistent with our observations of pairwise SNP differences (Fig. 4a),

202    we found that *C. difficile* had the lowest evolutionary rate, VRE and MRSA had intermediate

203    rates, and *P. aeruginosa* had the highest rate. Within each species group, however, we

204    observed a range of nucleotide substitution rates between the different STs that were sampled.

205    Rates overall varied nearly 100-fold among the species and STs we examined, from a minimum

206    of 0.40 SNPs/genome/year for *C. difficile* ST42, to 28.80 SNPs/genome/year for *P. aeruginosa*

207    ST179 (Fig. 4b, Supplementary Table 7). To understand how recombination might influence

208    these calculations, we used ClonalFrameML[24] to quantify the number of recombination events

209    per point mutation (R/Theta) for each ST across all species for which at least 10 different

210    isolates belonging to the same ST were sampled (Fig. 4c). MRSA genomes were found to have

211    the lowest rates of recombination, while *K. pneumoniae*, *E. coli*, and *A. baumannii* appeared to

212    have the highest rates. These data show that rates of nucleotide substitution and recombination

213    are variable across STs as well as across species; this variability should be considered when

214    assessing genomic similarity between isolates during epidemiologic investigations.

215    **Systematic analysis of antimicrobial resistance (AMR) genes uncovers broad and**

216    **species-specific trends**

217    AMR threatens the effective treatment and prevention of bacterial infections. To understand the

218    diversity and distribution of AMR genes among the 3,004 isolates we sampled, we identified

219    resistance genes within each genome by querying the ResFinder database with BLASTn[25]

220    (Supplementary Figure 4, Supplementary Table 8). The total number of AMR genes identified

221    per genome ranged from 0-19, with an average of 4.6 AMR genes per genome. The species

222    groups carrying the most AMR genes were *Klebsiella* spp. (average 13.1 AMR genes per

223    genome), *E. coli* (7.7 AMR genes per genome), and VRE (average 7.4 AMR genes per

224    genome) (Supplementary Table 8). We also classified each AMR gene by drug class, and

225    examined the distribution of AMR genes found in more than one species group (Fig. 5a).

226 Several genes encoding aminoglycoside and sulfonamide resistance were observed in the

227 majority of different species groups, suggesting that AMR genes for these antibiotic classes are

228 relatively widespread among bacterial pathogens within our hospital. The Gram-positive species

229 we collected (*C. difficile*, VRE, and MRSA) carried different AMR genes compared to the

230 sampled Gram-negative species, and all Gram-positive species were found to carry the

231 aminoglycoside resistance genes *aac(6')-aph(2')* and *aph(3')-III* and the tetracycline resistance

232 gene *tet(M)*, albeit at varying frequencies (Fig. 5a).

233 We next examined the co-occurrence of pairs of AMR genes across different species

234 groups (Fig. 5b). We found that the aminoglycoside resistance genes *aph(3")-Ib* and *aph(6)-Id*

235 were almost always found together, and co-occurred in eight different species groups (all Gram-

236 negative species groups except for *Burkholderia* spp., *Providencia* spp., and *Stenotrophomonas*

237 spp.). Both of these genes also frequently co-occurred with the sulfonamide resistance gene

238 *sul2* (Fig. 5b). A separate aminoglycoside resistance gene, *aac(6')-Ib-cr*, was found to

239 frequently co-occur with the narrow-spectrum beta-lactamase $bla_{OXA-1}$ as well as with the

240 extended-spectrum beta-lactamase (ESBL) $bla_{CTX-M-15}$. Finally, we examined the distribution of

241 ESBL and carbapenemase enzymes among the ESBL-producing *E. coli* and *Klebsiella* spp.

242 isolates that we sampled (Fig. 5c). The most frequently observed ESBL enzyme was CTX-M-15,

243 which was found in roughly half of all *E. coli* and *Klebsiella* spp. genomes (Fig. 5c). The other

244 half of isolates within each species group carried largely different enzymes from one another,

245 with most *E. coli* isolates carrying other CTX-M-type and a small number of TEM-type ESBLs,

246 while *Klebsiella* spp. isolates carried CTX-M-14 and SHV-type ESBLs. The carbapenemases

247 KPC-2, KPC-3, KPC-8, and KPC-31 were found almost entirely among *Klebsiella* spp. genomes

248 (Fig. 5c). These data highlight the abundant diversity of AMR genes carried by the bacteria in

249 our hospital, and can be useful for developing tailored treatment and prevention approaches for

250 different bacterial pathogens.

251 **Mobile genetic element (MGE) distribution and cargo**

252    MGEs are frequently found within the genomes of bacteria residing in the hospital environment,

253    and they often encode useful functions like AMR and virulence factors[26]. To assess the

254    presence of MGEs in our dataset in a systematic and unbiased manner, we used a previously

255    developed approach to identify nucleotide sequences with high homology (>99.9% identity over

256    at least 10Kb) that were present in genomes of different genomospecies[27] (Fig. 6a). This

257    approach resulted in the identification of 186 clusters of shared sequences, which were present

258    in 805 (26.8%) of the genomes in our dataset (Fig. 6b). While each of the 14 different species

259    groups we sampled contained at least one genome encoding a shared sequence, species

260    groups that were particularly enriched for shared sequences included *Klebsiella* spp., *P.*

261    *aeruginosa*, and *Stenotrophomonas* spp. (Fig. 6b). We next used comparisons with available

262    MGE databases and manual curation to assign an MGE type to each of the 186 clustered

263    sequences based on sequence homology to previously described MGEs (Fig. 6c). We identified

264    similar numbers of sequences that resembled insertion sequences (ISs) or transposons and that

265    resembled prophages or integrative conjugative elements (ICEs). Slightly more sequences

266    showed homology to plasmid sequences, and a large number of sequences resembled multiple

267    MGE types (Fig. 6c). Importantly, 53 (28.5%) shared sequence clusters could not be assigned

268    to an MGE type. Some of these sequences are likely fragments of larger MGEs that lacked

269    genetic features that would enable their classification. Alternately, some of these may constitute

270    novel MGEs.

271        To understand more about the cargo encoded by the putative MGEs we identified, we

272    first assessed the distribution of AMR genes among the 186 shared sequence clusters we

273    studied (Fig. 6d and Supplementary Table 9). Only 10/186 shared sequence clusters (5.4%)

274    carried AMR genes, however these clusters were found among 116/805 isolates (14.4%). The

275    most frequently observed AMR gene classes (which were each only present in four shared

276    sequence clusters) were sulfonamide and trimethoprim resistance, while aminoglycoside

277    resistance genes, tetracycline resistance genes, and beta-lactamases were each found in three

278  shared sequence clusters (Fig. 6d). We next examined the distribution of clusters of orthologous

279  groups of proteins (COG) categories among all genes present in all shared sequence clusters in

280  our dataset. A total of 938 genes (12.1% of all shared sequence cluster genes) had COG

281  categories assigned, and among these genes the two COG categories observed most

282  frequently were genes involved in replication, recombination and repair, and genes involved in

283  inorganic ion transport and metabolism (Fig. 6e). These data suggest that prominent cargo

284  among the shared sequences we identified included genes for MGE maintenance and

285  transmission, as well as genes required for the utilization of and resistance to heavy metals,

286  which pathogens frequently encounter in the hospital environment[28].

287

288  **Discussion**

289  HAIs place a large burden on healthcare systems by increasing patient morbidity, mortality, and

290  the cost of medical care. The broader aim of the EDS-HAT project is to improve the detection of

291  bacterial outbreaks in hospitals, and the project has been successful in this regard[14,15,22]. The

292  EDS-HAT project has also provided a large dataset of microbial genomes sampled from

293  thousands of patients within a single medical center over time. Here we highlight the genetic

294  diversity among bacterial pathogens causing HAIs; understanding this diversity can better

295  inform genomic epidemiology and outbreak investigations. As bacterial WGS becomes

296  increasingly routine in healthcare settings, this study also provides a baseline for future

297  comparisons, both at our center and elsewhere.

298      Using comparative genomics methods, we revealed the vast diversity among bacterial

299  pathogens within our hospital. We identified bacteria belonging to 97 different species, which

300  spanned 14 different species groups. We also identified 23 species which have not been

301  previously described, including potentially novel species of *Acinetobacter, Citrobacter, Proteus,*

302  *Providencia, Pseudomonas, Serratia* and *Stenotrophomonas.* A total of 41 isolates (1.4% of

303  sampled isolates) belonged to these novel species, which was a lower proportion than that

304    observed in a prior study of HAIs among ICU patients conducted in 2015[16]. This could be due to

305    additional species having been described in recent years, as well as different inclusion criteria

306    and study populations between the prior study and our own. Further investigation into these new

307    species can aid in the clinical diagnosis of bacteria causing infections.

308         Our finding that both ANI and PCA-A are effective at distinguishing between different

309    groups at both the species and sub-species levels is consistent with prior studies[29,30]. The 15 *P.*

310    *aeruginosa* isolates we identified as having 93-94% ANI with the remaining *P. aeruginosa*

311    population is also consistent with prior reports of the *P. aeruginosa* population[31]. Conversely, *S.*

312    *marcescens* is known to have a population structure comprised of multiple clades[32,33], however

313    we found that pairwise comparisons between some of these clades had less than 95% ANI,

314    suggesting a large degree of divergence and possible ongoing sub-speciation. We were also

315    able to use accessory gene content differences to distinguish between the dominant genetic

316    lineages of *C. difficile*, VRE *faecium*, and MRSA. Further investigation of these accessory genes

317    would likely enhance our understanding of how different genetic lineages are able to co-exist in

318    the same hospital, and could provide useful biomarkers for tracking lineages of interest.

319         Comparing within-patient versus between-patient genetic diversity can provide important

320    guidance in defining SNP cut-offs for outbreak investigations. We found that the number of

321    SNPs among genomes isolated from the same patient at different time points varied by species,

322    with within-patient SNPs being lowest for *C. difficile*, moderate for MRSA and VRE, and greatest

323    for *P. aeruginosa*. Differences between species likely reflect both genome size as well as the

324    biology of these organisms; for example, *C. difficile* can spend long periods of time in a non-

325    replicative spore state, while *P. aeruginosa* genomes are more than double the size of MRSA

326    and VRE genomes. The SNP distances among same-patient isolates we observed are

327    comparable to those used in outbreak investigations in our setting and elsewhere[14,34,35]. These

328    data demonstrate that same-patient genome pairs can be used to empirically determine genetic

329    similarity thresholds for genomic epidemiology purposes. Evolutionary rates assessed for the

330    four most common species in our hospital were also consistent with previous studies[36,37]. The

331    large variability in evolutionary rates between different species, however, further suggests that

332    different SNP cut-offs should be considered for different bacterial species for the purposes of

333    hospital outbreak investigations.

334         This study establishes the diversity of antimicrobial resistance genes among pathogenic

335    bacteria circulating at our hospital, and provides a point of comparison with other studies of

336    antibiotic resistance spread in the hospital environment[22,27,38,39]. We found that aminoglycoside

337    and sulfonamide resistance genes were highly abundant, and were found in the majority of

338    species that we sampled. Although the presence of aminoglycoside resistance is well

339    documented among both Gram-positive and Gram-negative bacteria—and more specifically

340    among the ESKAPE pathogens—less attention has been focused on sulfonamide resistance[40-

341    42]. The co-occurrence of *aph(3")-Ib, aph(6)-Id,* and *sul2* has been previously observed in a

342    variety of different genetic contexts, including in plasmids, integrative conjugative elements, and

343    chromosomal genomic islands[41,43]. Additionally, we found that the ESBL enzyme $bla_{CTX-M-15}$ was

344    widely distributed among both *E. coli* and *Klebsiella* spp. isolates, which is consistent with prior

345    reports[44]. Among the other ESBL-producing *E. coli* and *Klebsiella* spp. isolates collected, ESBL

346    enzymes were largely restricted to one species group or the other. Finally, while we did not

347    explicitly collect carbapenemase-producing organisms during this study period, a subset of the

348    ESBL-producing *E. coli* and *Klebsiella* spp. isolates collected also carried carbapenemase

349    enzymes. Co-occurrence of ESBL enzymes and carbapenemases was more frequent among

350    *Klebsiella* spp., especially ST258 *K. pneumoniae*[22].

351         This study also offers an overview of highly similar sequences (which we suspect largely

352    belong to MGEs) shared among the genomes of distantly related bacteria sampled from

353    patients residing in the same hospital environment. We found that *Enterobacteriaceae* such as

354    *Klebsiella* spp. and *Citrobacter* spp., as well as *P. aeruginosa* and *Stenotrophomonas* spp.,

355    were overrepresented among shared sequence clusters compared to their overall distribution in

356    the dataset. Most of the shared sequences identified in *Enterobacteriaceae* genomes resembled

357    sequences carried on plasmids, consistent with the frequent plasmid exchange known to

358    happen among species in this family[45]. On the other hand, many of the shared sequences

359    identified among *P. aeruginosa* and *Stenotrophomonas* spp. resembled prophages and

360    integrated conjugative elements, suggesting that these organisms may rely on different MGEs

361    to exchange genetic material. Somewhat surprisingly, our analysis identified fewer shared

362    sequences carrying AMR genes compared to a prior study we conducted within the same

363    hospital[27]. This may be due to our use of a longer sequence length cut-off for shared sequence

364    identification in this study, as AMR genes are known to be carried on smaller MGE units that

365    can rapidly shuffle, interchange, and mutate[46]. Finally, we found it notable that genes encoding

366    metal transport and resistance were frequently observed within the shared sequences we

367    identified. Inorganic ions are required for catalysis of many bacterial enzymes[47], and heavy

368    metals such as silver, copper, and mercury have long been used as disinfectants in hospitals[48].

369    Further study of MGEs encoding metal-interacting genes will be a focus of our future work.

370        This study had several limitations. The organisms we collected were pre-specified, and

371    certain groups, such as *Enterobacter* spp. or carbapenemase-producing organisms without a

372    noted ESBL phenotype, were not collected. Furthermore, our definition of "hospital-acquired

373    infections" was quite broad; some of the collected isolates likely represent commensal

374    organisms or pathogen colonization, rather than true infection. We also cannot say for sure

375    whether the sampled bacteria were acquired from the healthcare setting or not, as we only

376    considered bacterial isolates from clinical specimens and did not include environmental

377    sampling. Additionally, our 25-month collection window was quite short, thus we were unable to

378    draw conclusions regarding trends over time. Finally, the inclusion of both broad species groups

379    as well as more defined sets of specific pathogens made it difficult to conduct systematic

380    analyses or draw broader conclusions across the entire dataset. Nonetheless, the large number

381    of isolates collected offers a high-resolution view of the genomic diversity and evolution of

382   important bacterial pathogens found within our hospital. Our future work will include following

383   these bacterial populations over time, and comparing our results with similar studies conducted

384   in other settings.

385          In assessing the genomes of major infection-associated bacterial species isolated from

386   patients at our hospital, we have provided a longitudinal survey of the genomic diversity of

387   bacterial HAIs at a single clinical center. Our findings demonstrate that studying population

388   dynamics and evolution of these pathogens can inform genomics-based outbreak

389   investigations. In addition to forming a basis for future comparisons, this study also provides a

390   deeper understanding of the breadth of different species that cause HAIs, and demonstrates the

391   utility of systematic genome sequencing and comparative genomics analysis of clinical bacterial

392   isolates from hospitalized patients.

393

394   **Methods**

395   **Isolate collection**

396   Bacterial isolates were collected from the University of Pittsburgh Medical Center (UPMC)

397   Presbyterian Hospital, an adult tertiary care hospital with over 750 beds, 150 critical care unit

398   beds, more than 32,000 yearly inpatient admissions, and over 400 solid organ transplants per

399   year. Isolates were collected from November 2016 through November 2018 from admitted

400   patients as part of a prospective genomic epidemiology surveillance project called Enhanced

401   Detection System for Healthcare-Associated Transmission (EDS-HAT). Inclusion criteria were

402   hospital admission greater than two days before the culture date, and/or a recent inpatient or

403   outpatient UPMC hospital encounter in the 30 days before the culture date. A total of 3,004

404   isolates were included in this study (Table S1). The EDS-HAT project collected all organisms

405   meeting the above inclusion criteria and belonging to the following genera: *Acinetobacter* spp.*,*

406   *Burkholderia* spp., *Citrobacter* spp., *Proteus* spp., *Providencia* spp., *Pseudomonas* spp. *Serratia*

407   spp., and *Stenotrophomonas* spp. Isolate collection was limited to only toxin-producing strains

408    of *Clostridioides difficile*, vancomycin-resistant *Enterococcus* spp. (VRE), extended-spectrum

409    beta-lactamase (ESBL)-producing *Escherichia coli* and *Klebsiella* spp., and methicillin-resistant

410    *Staphylococcus aureus* (MRSA). This study was approved by the University of Pittsburgh

411    Institutional Review Board and was classified as being exempt from patient-informed consent.

412    **Whole genome sequencing and genome assembly**

413    Genomic DNA was extracted from pure overnight cultures of single bacterial colonies using a

414    Qiagen DNeasy Tissue Kit according to the manufacturer's instructions (Qiagen, Germantown,

415    MD). Illumina library construction and sequencing were conducted using an Illumina Nextera

416    DNA Sample Prep Kit with 150bp paired-end reads, and libraries were sequenced on the

417    NextSeq 550 sequencing platform (Illumina, San Diego, CA). Selected isolates were re-

418    sequenced with long-read technology on a MinION device (Oxford Nanopore Technologies,

419    Oxford, United Kingdom). Long-read sequencing libraries were prepared and multiplexed using

420    a rapid multiplex barcoding kit (catalog SQK-RBK004) and were sequenced on R9.4.1 flow

421    cells. Base-calling on raw reads was performed using Albacore v2.3.3 or Guppy v2.3.1 (Oxford

422    Nanopore Technologies, Oxford, UK).

423    Genome sequence analyses were performed on a BioLinux v8 server[49] using publicly

424    available genomic analysis tools wrapped together into a high-throughput genome analysis

425    pipeline. Briefly, Illumina sequencing data were processed with Trim Galore v0.6.1

426    (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)    to    remove    sequencing

427    adaptors, low-quality bases, and poor-quality reads. Kraken v1[50] taxonomic sequence

428    classification of raw reads was used to confirm species designation, and to rule out

429    contamination. Illumina reads were assembled with SPAdes v3.11[51]. Long-read sequence data

430    generated for other studies[22,27,39] were combined with Illumina data for the same isolate, and

431    hybrid assembly was conducted using unicycler v0.4.7 or v0.4.8-beta[52]. Assembled genomes

432    were annotated using Prokka v1.14 and assembly quality was verified using QUAST[53].

433    Genomes were included in the study if they had at least 35-fold Illumina read coverage, had

434    assemblies with ≤ 350 contigs, and had total genome lengths ± 25% of the median of all isolates

435    within each species group. Antimicrobial resistance and toxin genes were confirmed using

436    BLASTn in line with EDS-HAT study phenotypic inclusion criteria. Specifically, all *S. aureus*

437    genomes were confirmed to encode the *mecA* gene, all *E. faecalis* and *E. faecium* genomes

438    were confirmed to encode a VanA or VanB operon, all *E. coli* and *Klebsiella* spp. genomes were

439    confirmed to encode an identifiable extended-spectrum beta-lactamase (ESBL) enzyme, and all

440    *C. difficile* genomes were confirmed to encode either toxin A and/or toxin B genes.

441    **Classification of genomospecies and lineages**

442    Within each species group, genome assemblies from this study and reference genome

443    assemblies downloaded from the NCBI RefSeq database underwent pairwise average

444    nucleotide identity (ANI) analysis using FastANI v1.3[18]. Genomes with ANI values >95% then

445    underwent single-linkage hierarchical clustering using the hclust function from the R package

446    stats v3.6. Each ANI cluster was manually assessed and assigned to a species based on the

447    predominant nomenclature of genomes of type/reference strains within each cluster. Clusters

448    that did not contain reference genomes, or where reference genomes were only named at the

449    genus level, were named "genomospecies." Sequential numbers were appended to each

450    uncharacterized genomospecies within a species group. Species identified using ANI and

451    having greater than 100 isolates were further sub-divided into clades and lineages based on

452    multi-locus sequence typing (ST), or phylogenetic analysis. STs were determined from

453    assembled contigs using mlst v2 (https://github.com/tseemann/mlst). Species without a defined

454    ST scheme (*P. mirabilis* and *S. marcescens*) were classified into clades or lineages by grouping

455    isolates that shared <1000 core genome single nucleotide polymorphism (SNP) differences into

456    the same lineage, with SNPs identified using snippy (https://github.com/tseemann/snippy).

457    *Stenotrophomonas* genomospecies were named according to Gröschel et al.[54].

458    **Gene content and pangenome analyses**

459   Gene content matrices were obtained for all species groups with more than 50 isolates using

460   the pangenome analysis program roary v3.11[55]. Roary was run using a protein identity cut-off of

461   80% for genera containing multiple species, and a cut-off of 95% for individual species.

462   Pangenome collector's curves were generated for each species group by calculating the

463   number of unique genes present at increasing numbers of sampled genomes, with 1000

464   iterations of each sample size up to 250. Genetic clustering of genomes within species groups

465   based on variable gene content was calculated and visualized using principal component

466   analysis of accessory genes (PCA-A) using the R packages prcomp, vegan, and ggbiplot, with

467   matrices of gene presence/absence used as input. Genes that were present in all isolates,

468   present in only one isolate, or absent in only one isolate, were removed from analysis. PCA-A

469   coordinate plots were visualized using GraphPad Prism version 7.0c.

470   **Core genome SNP comparisons, phylogenetic trees, evolutionary rate and recombination**

471   **analyses**

472   Within each genus, species, ST, or clade, SNPs were identified using snippy

473   (https://github.com/tseemann/snippy). The most complete genome assembly (i.e. highest N50)

474   was used as a reference genome for SNP analysis. Core genome SNPs, defined as SNPs at

475   nucleotide positions shared across all genomes in the sample group being compared, were

476   used to calculate pairwise SNP distances and to generate maximum likelihood phylogenetic

477   trees. Trees were generated with RAxML v8.2 using the general time reversible model of

478   evolution (GTRCAT), Lewis correction for ascertainment bias, and 100 bootstrap replicates[56].

479   Unless otherwise specified, reported SNP distances refer to core genome SNPs for all isolates

480   belonging to the same ST. Pairwise SNP distances were visualized using the R package

481   ggplot2. Recombination and evolutionary rates were calculated for STs in four species groups

482   (*P. aeruginosa*, *Clostridioides difficile*, VRE and MRSA), and for STs within each group with

483   more than 25 isolates. Estimates of relative recombination rates (R/Theta) and average size of

484   recombinant sequences (delta) were assessed from core genome alignments using

485    ClonalFrameML v1.12[24] with default settings. The relative rate of recombination, which reflects

486    the number of nucleotide changes introduced by recombination relative to each point mutation

487    (r/m) was calculated as r/m = (R/Theta) × delta × v[24], where v is the average distance between

488    recombined sequences. A core genome alignment and recombination-corrected phylogenetic

489    tree were used to estimate evolutionary rates using TreeTime[23]. Isolates that were found to be

490    highly divergent from other isolates of the same ST (as revealed by an excess number of SNPs

491    separating them from other isolates) were removed from the analysis.

**Antibiotic resistance gene detection and analysis**

493    Acquired antimicrobial resistance genes were detected by querying genome assemblies against

494    the ResFinder database using BLASTn[25]. A gene was considered present if the BLASTn

495    percent identity multiplied by the sequence coverage was >80%. Resistance gene presence

496    was mapped to a global phylogenetic tree constructed from amino acid sequences of 120

497    ubiquitous protein coding genes from the Genome Taxonomy Database Tool Kit[57]. Resistance

498    gene co-occurrence was calculated using the %*% operator in R. This operator works by

499    identifying the cross-products between any two genes found in a matrix of resistance genes

500    identified in all isolates. The results were used to construct a relative frequency plot using the

501    ggplot2 package in R. To include only the most frequently co-occurring gene pairs in the plot, a

502    relative frequency of 80% and a combined frequency of 50% were used as cut-off thresholds.

503    Additionally, genes found in >250 isolates were excluded as they were suspected of not being

504    acquired resistance genes. ESBL and carbapenemase enzyme distributions were determined

505    by assigning enzyme types based on protein sequence, and only 100% protein sequence

506    matches are reported.

**Shared sequence detection and analysis**

508    Putative mobile genetic elements were identified by searching for sequences >10kb that were

509    present at high identity (>99.9%) in the genomes of isolates belonging to different species

510    (<95% ANI) using nucmer[58]. Sequences were organized into clusters using all-by-all BLASTn

511 v2.7.1[59], and clusters were visualized with Cytoscape v3.8.2[60]. Clustered shared sequences

512 were determined as resembling plasmids, insertion sequences (ISs), transposons, prophages,

513 or integrative conjugative elements by BLAST against complete plasmids from NCBI

514 databases[61], MobileElementFinder[62], PHASTER[63], ProphET[64] and ICEberg[65], as well as

515 comparison to the NCBI nr database and manual curation. Antimicrobial resistance genes in

516 clustered sequences were identified by BLASTn against the ResFinder database[25]. Clusters of

517 orthologous groups of proteins (COG) categories were assigned to genes present in one or

518 more clustered sequences, and the distribution of genes in each COG category was visualized

519 with the pie function in R.

520

521 **Data availability**

522 Raw sequencing reads and genome assemblies were submitted to the NCBI Sequence Read

523 Archive (SRA) and GenBank, with accession numbers listed in Table S1.

524

525 **References**

526 1  Magill, S. S. *et al.* Changes in Prevalence of Health Care-Associated Infections in U.S.
527    Hospitals. *N Engl J Med* **379**, 1732-1744, doi:10.1056/NEJMoa1801550 (2018).
528 2  Stone, P. W. Economic burden of healthcare-associated infections: an American
529    perspective. *Expert Rev Pharmacoecon Outcomes Res* **9**, 417-422,
530    doi:10.1586/erp.09.53 (2009).
531 3  Centers for Disease Control and Prevention. *Current HAI Progress Report*,
532    <https://www.cdc.gov/hai/data/portal/progress-report.html> (2020, Dec 02).
533 4  Rice, L. B. Federal funding for the study of antimicrobial resistance in nosocomial
534    pathogens: no ESKAPE. *J Infect Dis* **197**, 1079-1081, doi:10.1086/533452 (2008).
535 5  Centers for Disease Control and Prevention. ANTIBIOTIC RESISTANCE THREATS in
536    the United States, 2013. (CDC, 2013).
537 6  Lax, S. *et al.* Bacterial colonization and succession in a newly opened hospital. *Sci
538    Transl Med* **9**, doi:10.1126/scitranslmed.aah6500 (2017).
539 7  Curry, S. R. *et al.* Use of multilocus variable number of tandem repeats analysis
540    genotyping to determine the role of asymptomatic carriers in Clostridium difficile
541    transmission. *Clin Infect Dis* **57**, 1094-1102, doi:10.1093/cid/cit475 (2013).
542 8  Kanamori, H., Rutala, W. A. & Weber, D. J. The Role of Patient Care Items as a Fomite
543    in Healthcare-Associated Outbreaks and Infection Prevention. *Clin Infect Dis* **65**, 1412-
544    1419, doi:10.1093/cid/cix462 (2017).
545 9  Santajit, S. & Indrawattana, N. Mechanisms of Antimicrobial Resistance in ESKAPE
546    Pathogens. *Biomed Res Int* **2016**, 2475067, doi:10.1155/2016/2475067 (2016).

547  10  Quainoo, S. *et al.* Whole-Genome Sequencing of Bacterial Pathogens: the Future of
548      Nosocomial Outbreak Analysis. *Clin Microbiol Rev* **30**, 1015-1063,
549      doi:10.1128/CMR.00016-17 (2017).
550  11  Peacock, S. J., Parkhill, J. & Brown, N. M. Changing the paradigm for hospital outbreak
551      detection by leading with genomic surveillance of nosocomial pathogens. *Microbiology*
552      *(Reading)* **164**, 1213-1219, doi:10.1099/mic.0.000700 (2018).
553  12  Sundermann, A. J. *et al.* Automated data mining of the electronic health record for
554      investigation of healthcare-associated outbreaks. *Infect Control Hosp Epidemiol* **40**, 314-
555      319, doi:10.1017/ice.2018.343 (2019).
556  13  Miller, J. K. *et al.* Statistical outbreak detection by joining medical records and pathogen
557      similarity. *J Biomed Inform* **91**, 103126, doi:10.1016/j.jbi.2019.103126 (2019).
558  14  Sundermann, A. J. *et al.* Outbreak of Pseudomonas aeruginosa Infections from a
559      Contaminated Gastroscope Detected by Whole Genome Sequencing Surveillance. *Clin*
560      *Infect Dis*, doi:10.1093/cid/ciaa1887 (2020).
561  15  Sundermann, A. J. *et al.* Outbreak of Vancomycin-resistant Enterococcus faecium in
562      Interventional Radiology: Detection Through Whole-genome Sequencing-based
563      Surveillance. *Clin Infect Dis* **70**, 2336-2343, doi:10.1093/cid/ciz666 (2020).
564  16  Roach, D. J. *et al.* A Year of Infection in the Intensive Care Unit: Prospective Whole
565      Genome Sequencing of Bacterial Clinical Isolates Reveals Cryptic Transmissions and
566      Novel Microbiota. *PLoS Genet* **11**, e1005413, doi:10.1371/journal.pgen.1005413 (2015).
567  17  Mosquera-Rendon, J. *et al.* Pangenome-wide and molecular evolution analyses of the
568      Pseudomonas aeruginosa species. *BMC Genomics* **17**, 45, doi:10.1186/s12864-016-
569      2364-4 (2016).
570  18  Jain, C., Rodriguez, R. L., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
571      throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.
572      *Nat Commun* **9**, 5114, doi:10.1038/s41467-018-07641-9 (2018).
573  19  Jeukens, J. *et al.* A Pan-Genomic Approach to Understand the Basis of Host Adaptation
574      in Achromobacter. *Genome Biol Evol* **9**, 1030-1046, doi:10.1093/gbe/evx061 (2017).
575  20  Potter, R. F., Burnham, C. D. & Dantas, G. In Silico Analysis of Gardnerella
576      Genomospecies Detected in the Setting of Bacterial Vaginosis. *Clin Chem* **65**, 1375-
577      1387, doi:10.1373/clinchem.2019.305474 (2019).
578  21  Roy, P. H. *et al.* Complete genome sequence of the multiresistant taxonomic outlier
579      Pseudomonas aeruginosa PA7. *PLoS One* **5**, e8842, doi:10.1371/journal.pone.0008842
580      (2010).
581  22  Marsh, J. W. *et al.* Evolution of Outbreak-Causing Carbapenem-Resistant Klebsiella
582      pneumoniae ST258 at a Tertiary Care Hospital over 8 Years. *mBio* **10**,
583      doi:10.1128/mBio.01945-19 (2019).
584  23  Sagulenko, P., Puller, V. & Neher, R. A. TreeTime: Maximum-likelihood phylodynamic
585      analysis. *Virus Evol* **4**, vex042, doi:10.1093/ve/vex042 (2018).
586  24  Didelot, X. & Wilson, D. J. ClonalFrameML: efficient inference of recombination in whole
587      bacterial genomes. *PLoS Comput Biol* **11**, e1004041, doi:10.1371/journal.pcbi.1004041
588      (2015).
589  25  Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *J Antimicrob*
590      *Chemother* **67**, 2640-2644, doi:10.1093/jac/dks261 (2012).
591  26  Lerminiaux, N. A. & Cameron, A. D. S. Horizontal transfer of antibiotic resistance genes
592      in clinical environments. *Can J Microbiol* **65**, 34-44, doi:10.1139/cjm-2018-0275 (2019).
593  27  Evans, D. R. *et al.* Systematic detection of horizontal gene transfer across genera
594      among multidrug-resistant bacteria in a single hospital. *Elife* **9**, doi:10.7554/eLife.53886
595      (2020).
596  28  McDonnell, G. & Russell, A. D. Antiseptics and disinfectants: activity, action, and
597      resistance. *Clin Microbiol Rev* **12**, 147-179 (1999).

598  29  McNally, A. *et al.* Combined Analysis of Variation in Core, Accessory and Regulatory
599      Genome Regions Provides a Super-Resolution View into the Evolution of Bacterial
600      Populations. *PLoS Genet* **12**, e1006280, doi:10.1371/journal.pgen.1006280 (2016).
601  30  Inglin, R. C., Meile, L. & Stevens, M. J. A. Clustering of Pan- and Core-genome of
602      Lactobacillus provides Novel Evolutionary Insights for Differentiation. *BMC Genomics*
603      **19**, 284, doi:10.1186/s12864-018-4601-5 (2018).
604  31  Freschi, L. *et al.* The Pseudomonas aeruginosa Pan-Genome Provides New Insights on
605      Its Population Structure, Horizontal Gene Transfer, and Pathogenicity. *Genome Biol Evol*
606      **11**, 109-120, doi:10.1093/gbe/evy259 (2019).
607  32  Moradigaravand, D., Boinett, C. J., Martin, V., Peacock, S. J. & Parkhill, J. Recent
608      independent emergence of multiple multidrug-resistant Serratia marcescens clones
609      within the United Kingdom and Ireland. *Genome Res* **26**, 1101-1109,
610      doi:10.1101/gr.205245.116 (2016).
611  33  Abreo, E. & Altier, N. Pangenome of Serratia marcescens strains from nosocomial and
612      environmental origins reveals different populations and the links between them. *Sci Rep*
613      **9**, 46, doi:10.1038/s41598-018-37118-0 (2019).
614  34  Eyre, D. W. *et al.* Diverse sources of C. difficile infection identified on whole-genome
615      sequencing. *N Engl J Med* **369**, 1195-1205, doi:10.1056/NEJMoa1216064 (2013).
616  35  Coll, F. *et al.* Definition of a genetic relatedness cutoff to exclude recent transmission of
617      meticillin-resistant Staphylococcus aureus: a genomic epidemiology analysis. *Lancet*
618      *Microbe* **1**, e328-e335, doi:10.1016/S2666-5247(20)30149-X (2020).
619  36  Miyoshi-Akiyama, T. *et al.* Emergence and Spread of Epidemic Multidrug-Resistant
620      Pseudomonas aeruginosa. *Genome Biol Evol* **9**, 3238-3245, doi:10.1093/gbe/evx243
621      (2017).
622  37  Didelot, X. *et al.* Microevolutionary analysis of Clostridium difficile genomes to
623      investigate transmission. *Genome Biol* **13**, R118, doi:10.1186/gb-2012-13-12-r118
624      (2012).
625  38  van Duin, D. *et al.* Molecular and clinical epidemiology of carbapenem-resistant
626      Enterobacterales in the USA (CRACKLE-2): a prospective cohort study. *Lancet Infect*
627      *Dis* **20**, 731-741, doi:10.1016/S1473-3099(19)30755-8 (2020).
628  39  Babiker, A. *et al.* Clinical and Genomic Epidemiology of Carbapenem-Nonsusceptible
629      Citrobacter spp. at a Tertiary Health Care Center over 2 Decades. *J Clin Microbiol* **58**,
630      doi:10.1128/JCM.00275-20 (2020).
631  40  World Health Organization. GLASS whole-genome sequencing for surveillance of
632      antimicrobial resistance. (22 September 2020).
633  41  Ramirez, M. S. & Tolmasky, M. E. Aminoglycoside modifying enzymes. *Drug Resist*
634      *Updat* **13**, 151-171, doi:10.1016/j.drup.2010.08.003 (2010).
635  42  De Oliveira, D. M. P. *et al.* Antimicrobial Resistance in ESKAPE Pathogens. *Clin*
636      *Microbiol Rev* **33**, doi:10.1128/CMR.00181-19 (2020).
637  43  Pal, C., Bengtsson-Palme, J., Kristiansson, E. & Larsson, D. G. Co-occurrence of
638      resistance genes to antibiotics, biocides and metals reveals novel insights into their co-
639      selection potential. *BMC Genomics* **16**, 964, doi:10.1186/s12864-015-2153-5 (2015).
640  44  Canton, R., Gonzalez-Alba, J. M. & Galan, J. C. CTX-M Enzymes: Origin and Diffusion.
641      *Front Microbiol* **3**, 110, doi:10.3389/fmicb.2012.00110 (2012).
642  45  Redondo-Salvo, S. *et al.* Pathways for horizontal gene transfer in bacteria revealed by a
643      global map of their plasmids. *Nat Commun* **11**, 3602, doi:10.1038/s41467-020-17278-2
644      (2020).
645  46  Partridge, S. R., Kwong, S. M., Firth, N. & Jensen, S. O. Mobile Genetic Elements
646      Associated with Antimicrobial Resistance. *Clin Microbiol Rev* **31**,
647      doi:10.1128/CMR.00088-17 (2018).

648  47   Chandrangsu, P., Rensing, C. & Helmann, J. D. Metal homeostasis and resistance in
649       bacteria. *Nat Rev Microbiol* **15**, 338-350, doi:10.1038/nrmicro.2017.15 (2017).
650  48   Villapun, V. M., Dover, L. G., Cross, A. & Gonzalez, S. Antibacterial Metallic Touch
651       Surfaces. *Materials (Basel)* **9**, doi:10.3390/ma9090736 (2016).
652  49   Field, D. *et al.* Open software for biologists: from famine to feast. *Nat Biotechnol* **24**, 801-
653       803, doi:10.1038/nbt0706-801 (2006).
654  50   Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification
655       using exact alignments. *Genome Biol* **15**, R46, doi:10.1186/gb-2014-15-3-r46 (2014).
656  51   Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to
657       single-cell sequencing. *J Comput Biol* **19**, 455-477, doi:10.1089/cmb.2012.0021 (2012).
658  52   Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial
659       genome assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**,
660       e1005595, doi:10.1371/journal.pcbi.1005595 (2017).
661  53   Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for
662       genome assemblies. *Bioinformatics* **29**, 1072-1075, doi:10.1093/bioinformatics/btt086
663       (2013).
664  54   Groschel, M. I. *et al.* The phylogenetic landscape and nosocomial spread of the
665       multidrug-resistant opportunist Stenotrophomonas maltophilia. *Nat Commun* **11**, 2044,
666       doi:10.1038/s41467-020-15123-0 (2020).
667  55   Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.
668       *Bioinformatics* **31**, 3691-3693, doi:10.1093/bioinformatics/btv421 (2015).
669  56   Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
670       large phylogenies. *Bioinformatics* **30**, 1312-1313, doi:10.1093/bioinformatics/btu033
671       (2014).
672  57   Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
673       substantially revises the tree of life. *Nat Biotechnol* **36**, 996-1004, doi:10.1038/nbt.4229
674       (2018).
675  58   Marcais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS
676       Comput Biol* **14**, e1005944, doi:10.1371/journal.pcbi.1005944 (2018).
677  59   Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
678       search tool. *J Mol Biol* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2 (1990).
679  60   Shannon, P. *et al.* Cytoscape: a software environment for integrated models of
680       biomolecular interaction networks. *Genome Res* **13**, 2498-2504, doi:10.1101/gr.1239303
681       (2003).
682  61   Che, Y. *et al.* Conjugative plasmids interact with insertion sequences to shape the
683       horizontal transfer of antimicrobial resistance genes. *Proc Natl Acad Sci U S A* **118**,
684       doi:10.1073/pnas.2008731118 (2021).
685  62   Johansson, M. H. K. *et al.* Detection of mobile genetic elements associated with
686       antibiotic resistance in Salmonella enterica using a newly developed web tool:
687       MobileElementFinder. *J Antimicrob Chemother* **76**, 101-109, doi:10.1093/jac/dkaa390
688       (2021).
689  63   Arndt, D. *et al.* PHASTER: a better, faster version of the PHAST phage search tool.
690       *Nucleic Acids Res* **44**, W16-21, doi:10.1093/nar/gkw387 (2016).
691  64   Reis-Cunha, J. L., Bartholomeu, D. C., Manson, A. L., Earl, A. M. & Cerqueira, G. C.
692       ProphET, prophage estimation tool: A stand-alone prophage sequence prediction tool
693       with  self-updating  reference  database.  *PLoS  One*  **14**,  e0223364,
694       doi:10.1371/journal.pone.0223364 (2019).
695  65   Bi, D. *et al.* ICEberg: a web-based resource for integrative and conjugative elements
696       found in Bacteria. *Nucleic Acids Res* **40**, D621-626, doi:10.1093/nar/gkr846 (2012).
697

708

709    **Figure Legends**

710    **Figure 1. Species and body site distribution of 3,004 clinical bacterial isolates from**

711    **hospitalized patients.** Isolates were collected from a single hospital over 25 months as part of

712    the Enhanced Detection System for Healthcare-Associated Transmission (EDS-HAT) project.

713    Pie charts show the distribution of isolates belonging to 14 different species groups collected

714    from different types of clinical specimens.

715    **Figure 2. Genome length and pangenome size among sampled species.** (A) Distribution of

716    genome lengths of isolates belonging to each species group, ordered from shortest to longest

717    median genome length. Vertical lines show median values. (B) Pangenome collection curves for

718    up to 250 genomes from genera containing multiple species and with at least 50 genomes

719    collected. Pangenomes were generated by Roary with an 80% protein identity cut-off. (C)

720    Pangenome collection curves for up to 250 genomes from species with at least 40 genomes

721    collected. Pangenomes were generated by Roary with an 95% protein identity cut-off. Curves

722    show the mean pan-genome size and shading shows the standard deviation.

723 **Figure 3. Average nucleotide identity (ANI) and principal component analysis of**

724 **accessory genes (PCA-A) distinguish between and within species.** (A) Phylogeny and

725 pairwise ANI values for *Citrobacter* spp. sampled by EDS-HAT. Grey shading indicates ANI

726 values >95%, with darker shading showing higher identity. (B) PCA-A plot for *Citrobacter*

727 species with >2 isolates. (C) Pairwise ANI distribution of *S. marcescens* isolate genomes,

728 showing pairwise ANI comparisons between isolates in different clades that fall below the

729 species cut-off (95% ANI, vertical dashed line). (D) PCA-A plot for *S. marcescens* isolates,

730 showing clear separation of five distinct clades. (E-G) PCA-A plots for dominant sequence types

731 (STs) of *C. difficile* (E), *E. faecium* (F), and *S. aureus* (G).

732 **Figure 4. Pairwise SNP distances and genome evolution vary between species.** (A)

733 Comparison of within-patient, within-cluster, and between-patient single nucleotide

734 polymorphisms (SNPs) for select species. Pairwise comparisons are shown for all isolate pairs

735 belonging to the same sequence type (ST) within each species. (B) Genome evolution rates for

736 the dominant STs within *C. difficile* (CD), vancomycin-resistant *E. faecium* (VRE), methicillin-

737 resistant *S. aureus* (MRSA) and *P. aeruginosa* (PSA). Isolates belonging to the four largest STs

738 (three largest for MRSA) of each species were considered, and nucleotide substitution rate

739 (SNPs/genome/year) was calculated for each ST separately. Individual data points are labeled

740 with the corresponding ST, and boxes show the median, $25^{th}$ and $75^{th}$ percentiles. (C)

741 Recombination events per mutation (R/Theta) for select species. Each data point represents a

742 distinct ST, and data are grouped by species. STs with at least 10 isolates are shown. Boxes

743 show the median, $25^{th}$ and $75^{th}$ percentiles. PRO=*P. mirabilis*, SER=*S. marcescens*, KLP=*K.*

744 *pneumoniae*, EC=*E. coli*, ACIN=*A. baumannii*.

745 **Figure 5. Antimicrobial resistance gene abundance and diversity.** (A) Prevalence of

746 resistance genes found in more than one species group. Genes are grouped by antibiotic class,

747 and grey shading shows the prevalence of each gene within and across each group. Darker

748 shading indicates higher prevalence. ACIN=*Acinetobacter* spp.; KL=*Klebsiella* spp.;

749  CB=*Citrobacter* spp.; EC=*E. coli*; PRV=*Providencia* spp.; PR=*Proteus* spp.; SER=*Serratia* spp.;

750  PSA=*P. aeruginosa*; PSB=*Pseudomonas* spp.; STEN=*Stenotrophomonas* spp.;

751  BC=*Burkholderia* spp.; VRE=vancomycin-resistant *Enterococcus* spp.; MRSA=methicillin-

752  resistant *S. aureus*; CD=*C. difficile*. (B) Resistance gene co-occurrence. Relative frequency

753  versus number of genomes is plotted for pairs of resistance genes that co-occur at ≥50%

754  relative frequency. Blue dots indicate AMR genes in the same drug class, while orange dots

755  indicate genes in different classes. The size of each dot corresponds to the number of different

756  species groups found to carry each pair. AMR gene pairs found in ≥4 different species groups

757  are labeled. (C) Distribution of extended-spectrum beta-lactamase (ESBL) and carbapenemase

758  enzymes among *E. coli* and *Klebsiella* spp. isolates.

759  **Figure 6. Mobile genetic element (MGE) distribution and cargo.** (A) Clusters of putative

760  MGEs identified in 3,004 study isolate genomes. Nodes within each cluster correspond to

761  bacterial isolates, and are color coded by species group (color key provided in panel B). (B)

762  Distribution of isolates in the entire dataset (left) versus isolates encoding one or more putative

763  MGEs (right). (C) Distribution of putative MGEs resembling plasmid, IS/transposon, or

764  prophage/ICE sequences, determined by nucleotide sequence comparisons and manual

765  curation. (D) Distribution of antimicrobial resistance (AMR) genes detected among 186 putative

766  MGEs. (E) Distribution of clusters of orthologous groups of proteins (COG) categories of MGE

767  genes with COG categories assigned.

768  **Figure S1. Average nucleotide identity (ANI) and principal components analysis of**

769  **accessory genes (PCA-A) among diverse species groups sampled by EDS-HAT.** (A)

770  Phylogenetic tree with pairwise ANI values and (B) PCA-A plot for *Acinetobacter* spp. (C)

771  Phylogeny and ANI of *Burkholderia* spp*.*, (D) *Providencia* spp., (E) *Pseudomonas* spp., and (F)

772  *Stenotrophomonas* spp. (G) PCA-A plot for *Stenotrophomonas* spp. Grey shading indicates ANI

773  values >95%, with darker shading showing higher identity. PCA-A plots include species with >2

774  isolates.

775 **Figure S2. Average nucleotide identity (ANI) comparisons of *P. aeruginosa* isolates.**

776 Histogram of pairwise ANI values for 863 *P. aeruginosa* isolate genomes sampled by EDS-HAT.

777 Dashed vertical line indicates 95% ANI. Comparisons in red are between isolates in *P.*

778 *aeruginosa* Groups 1 or 2 versus isolates in the PA7-like Group 3, which appear to belong to a

779 distinct genomospecies.

780 **Figure S3. Average nucleotide identity (ANI) comparisons of *S. marcescens* isolates.** (A)

781 Phylogeny and ANI of 177 *S. marcescens* isolates sampled by EDS-HAT. Grey shading

782 indicates ANI values >95%, with darker shading showing higher identity. White indicates ANI

783 values <95%. (B) Distribution of pairwise ANI values for *S. marcescens* isolates belonging to the

784 same or different clades, broken down into pairwise clade comparisons. All comparisons

785 between isolates in Clade A vs. Clade C and Clade A vs. Clade E fall below the standard

786 species cutoff of 95%.

787 **Figure S4. Distribution of antimicrobial resistance (AMR) genes among 3,004 clinical**

788 **bacterial isolates from hospitalized patients.** Resistance genes were identified by BLASTn

789 comparison to the ResFinder database. Isolates are ordered according to their phylogenetic

790 placement using the amino acid sequences of 120 ubiquitous protein-coding genes from the

791 Genome Taxonomy Database Tool Kit. "# Gene" shows the number of AMR genes per genome,

792 with darker shading indicating more AMR genes. The matrix shows the presence or absence of

793 202 AMR genes, grouped by antibiotic class. Heat maps at the top show the number of species

794 groups and total number of genomes encoding each gene, with darker shading indicating higher

795 numbers. Raw data used to make the matrix are available in Table S3.

**Fig. 1. Species and body site distribution of 3,004 clinical bacterial isolates from hospitalized patients.** Isolates were collected from a single hospital over 25 months as part of the Enhanced Detection System for Healthcare-Associated Transmission (EDS-HAT) project. Pie charts show the distribution of isolates belonging to 14 different species groups collected from different types of clinical specimens.
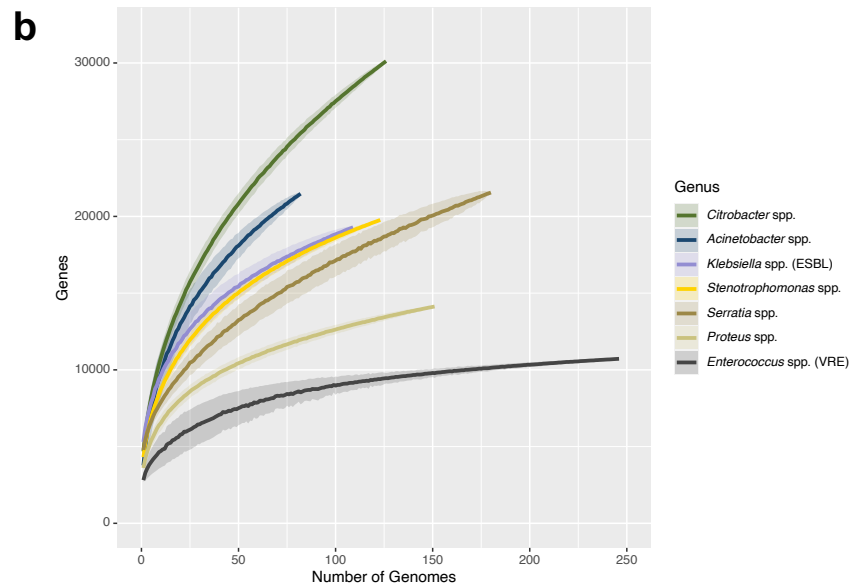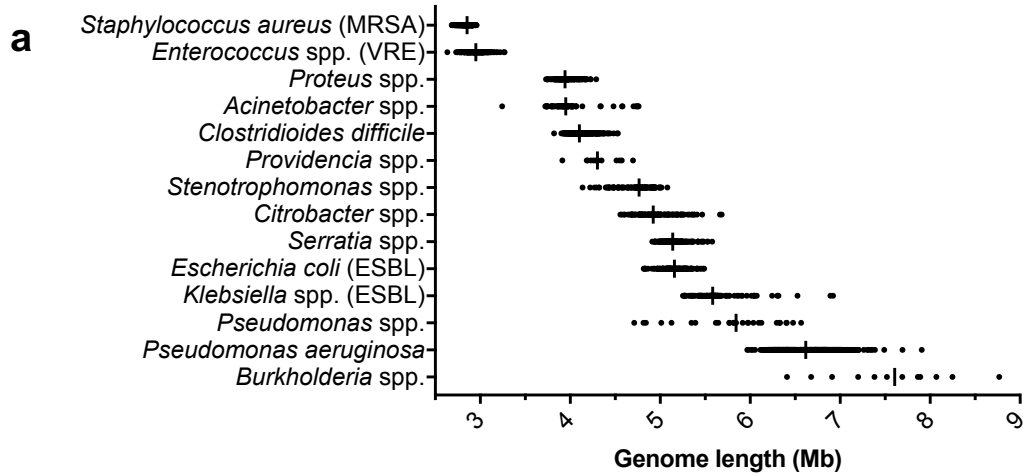
Fig. 2. Genome length and pangenome size among sampled species. a, Distribution of genome lengths of isolates belonging to each species group, ordered from shortest to longest median genome length. Vertical lines show median values. b, Pangenome collection curves for up to 250 genomes from genera containing multiple species and with at least 50 genomes collected. Pangenomes were generated by Roary with an 80% protein identity cut-off. c, Pangenome collection curves for up to 250 genomes from species with at least 40 genomes collected. Pangenomes were generated by Roary with an 95% protein identity cut-off. Curves show the mean pan-genome size and shading shows the standard deviation.

**Fig. 3. Average nucleotide identity (ANI) and principal component analysis of accessory genes (PCA-A) distinguish between and within species. a**, Phylogeny and pairwise ANI values for *Citrobacter* spp. sampled by EDS-HAT. Grey shading indicates ANI values >95%, with darker shading showing higher identity. **b**, PCA-A plot for *Citrobacter* species with >2 isolates. **c**, Pairwise ANI distribution of *S. marcescens* isolate genomes, showing pairwise ANI comparisons between isolates in different clades that fall below the species cut-off (95% ANI, vertical dashed line). **d**, PCA-A plot for *S. marcescens* isolates, showing clear separation of five distinct clades. **e-g**, PCA-A plots for dominant sequence types (STs) of *C. difficile* (**e**), *E. faecium* (**f**), and *S. aureus* (**g**).

**Fig. 4. Pairwise SNP distances and genome evolution vary between species. a**, Comparison of within-patient, within-cluster, and between-patient single nucleotide polymorphisms (SNPs) for select species. Pairwise comparisons are shown for all isolate pairs belonging to the same sequence type (ST) within each species. **b**, Genome evolution rates for the dominant STs within *C. difficile* (CD), vancomycin-resistant *E. faecium* (VRE), methicillin-resistant *S. aureus* (MRSA) and *P. aeruginosa* (PSA). Isolates belonging to the four largest STs (three largest for MRSA) of each species were considered, and nucleotide substitution rate (SNPs/genome/year) was calculated for each ST separately. Individual data points are labeled with the corresponding ST, and boxes show the median, 25th and 75th percentiles. **c**, Recombination events per mutation (R/Theta) for select species. Each data point represents a distinct ST, and data are grouped by species. STs with at least 10 isolates are shown. Boxes show the median, 25th and 75th percentiles. PRO=*P. mirabilis*, SER=*S. marcescens*, KLP=*K. pneumoniae*, EC=*E. coli*, ACIN=*A. baumannii*.

**Fig. 5. Antimicrobial resistance gene abundance and diversity. a**, Prevalence of resistance genes found in more than one species group. Genes are grouped by antibiotic class, and grey shading shows the prevalence of each gene within and across each group. Darker shading indicates higher prevalence. ACIN=*Acinetobacter* spp.; KL=*Klebsiella* spp.; CB=*Citrobacter* spp.; EC=*E. coli*; PRV=*Providencia* spp.; PR=*Proteus* spp.; SER=*Serratia* spp.; PSA=*P. aeruginosa*; PSB=*Pseudomonas* spp.; STEN=*Stenotrophomonas* spp.; BC=*Burkholderia* spp.; VRE=vancomycin-resistant *Enterococcus* spp.; MRSA=methicillin-resistant *S. aureus*; CD=*C. difficile*. **b**, Resistance gene co-occurrence. Relative frequency versus number of genomes is plotted for pairs of resistance genes that co-occur at ≥50% relative frequency. Blue dots indicate AMR genes in the same drug class, while orange dots indicate genes in different classes. The size of each dot corresponds to the number of different species groups found to carry each pair. AMR gene pairs found in ≥4 different species groups are labeled. **c**, Distribution of extended-spectrum beta-lactamase (ESBL) and carbapenemase enzymes among *E. coli* and *Klebsiella* spp. isolates.
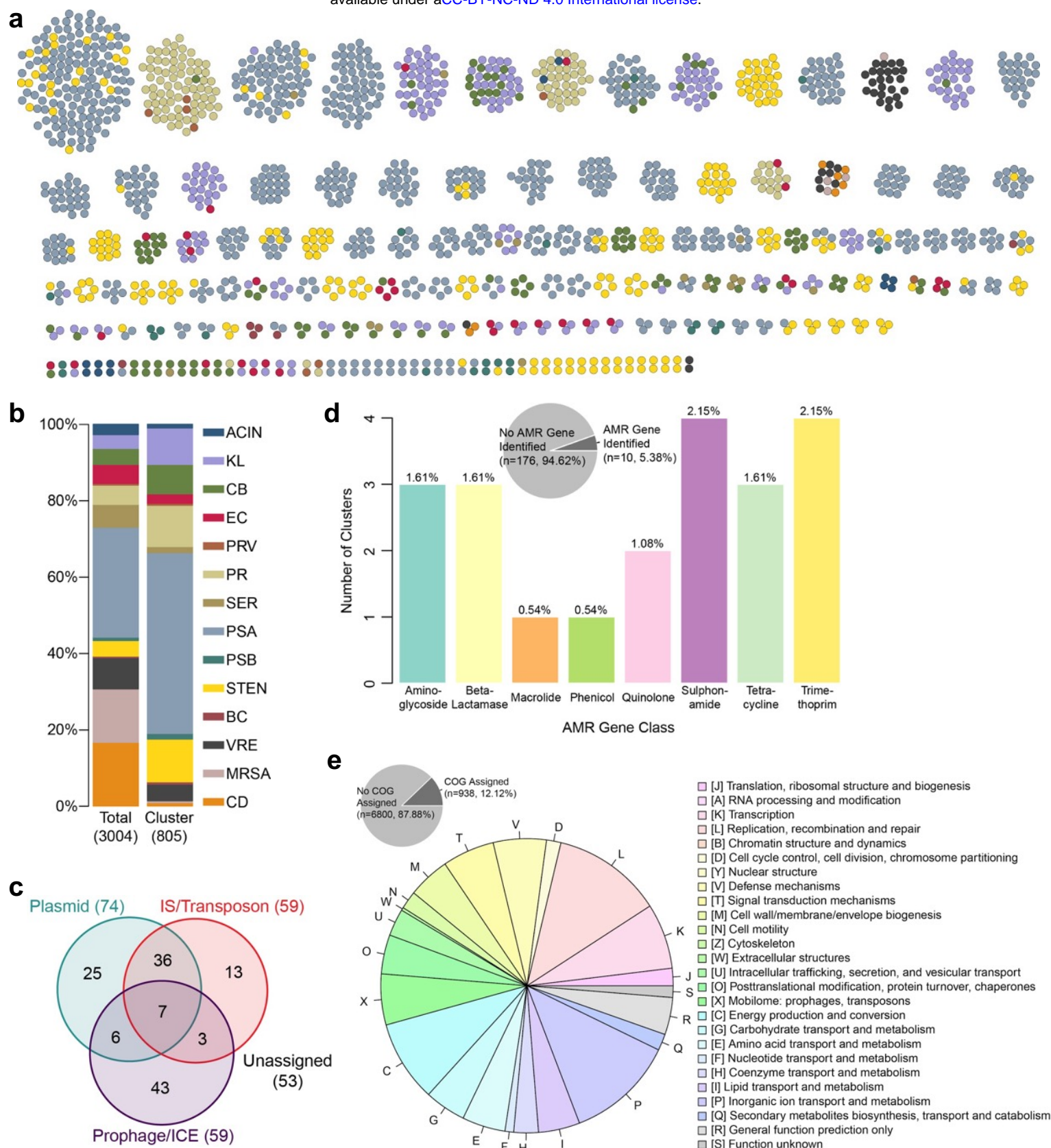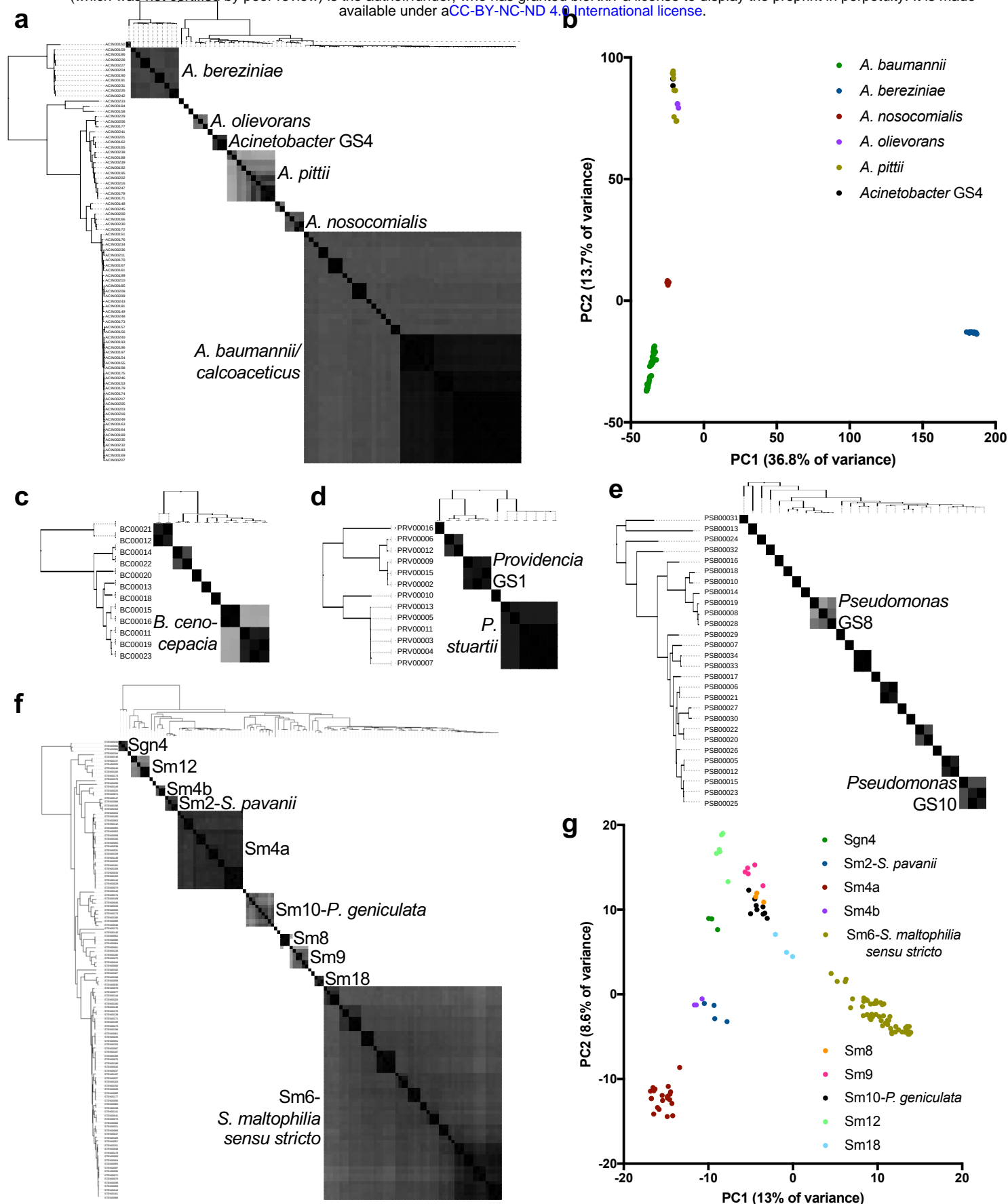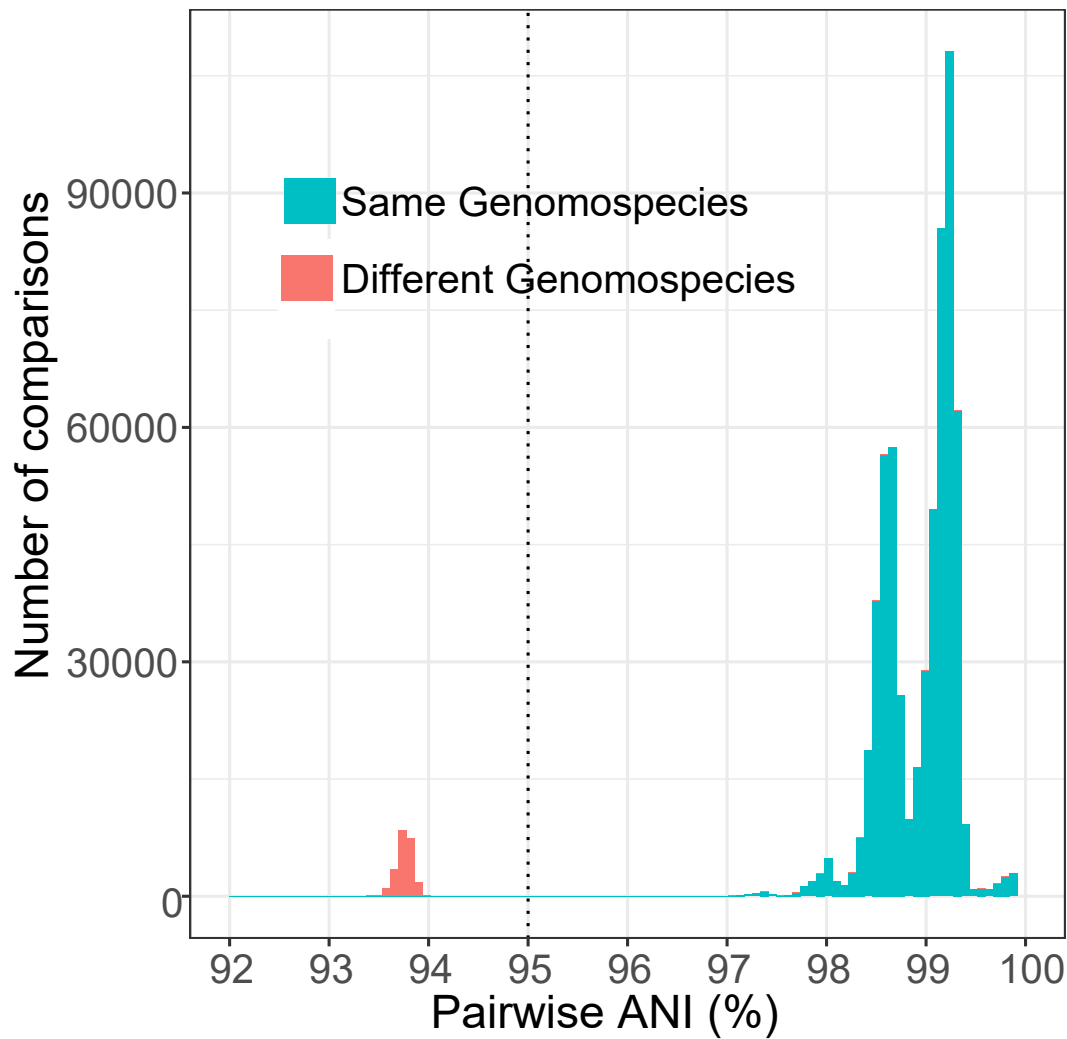
**Fig. 6. Mobile genetic element (MGE) distribution and cargo. a**, Clusters of putative MGEs identified in 3,004 study isolate genomes. Nodes within each cluster correspond to bacterial isolates, and are color coded by species group (color key provided in **b**). **b**, Distribution of isolates in the entire dataset (left) versus isolates encoding one or more putative MGEs (right). **c**, Distribution of putative MGEs resembling plasmid, IS/transposon, or prophage/ICE sequences, determined by nucleotide sequence comparisons and manual curation. **d**, Distribution of antimicrobial resistance (AMR) genes detected among 186 putative MGEs. **e**, Distribution of clusters of orthologous groups of proteins (COG) categories of MGE genes with COG categories assigned.
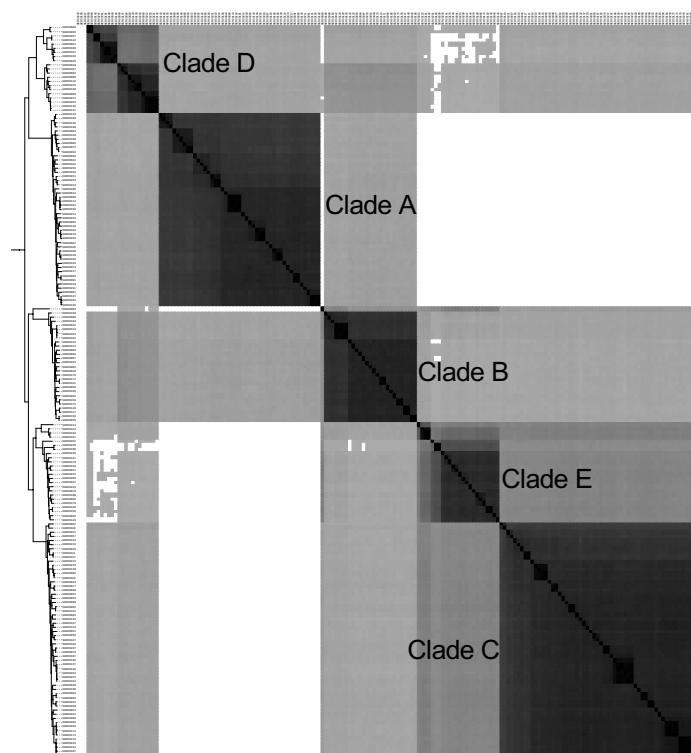
**Supplementary Fig. 1. Average nucleotide identity (ANI) and principal components analysis of accessory genes (PCA-A) among diverse species groups sampled by EDS-HAT. a**, Phylogenetic tree with pairwise ANI values and **b**, PCA-A plot for *Acinetobacter* spp. **c**, Phylogeny and ANI of *Burkholderia* spp., **d**, *Providencia* spp., **e**, *Pseudomonas* spp., and **f**, *Stenotrophomonas* spp. **g**, PCA-A plot for *Stenotrophomonas* spp. Grey shading indicates ANI values >95%, with darker shading showing higher identity. PCA-A plots include species with >2 isolates.
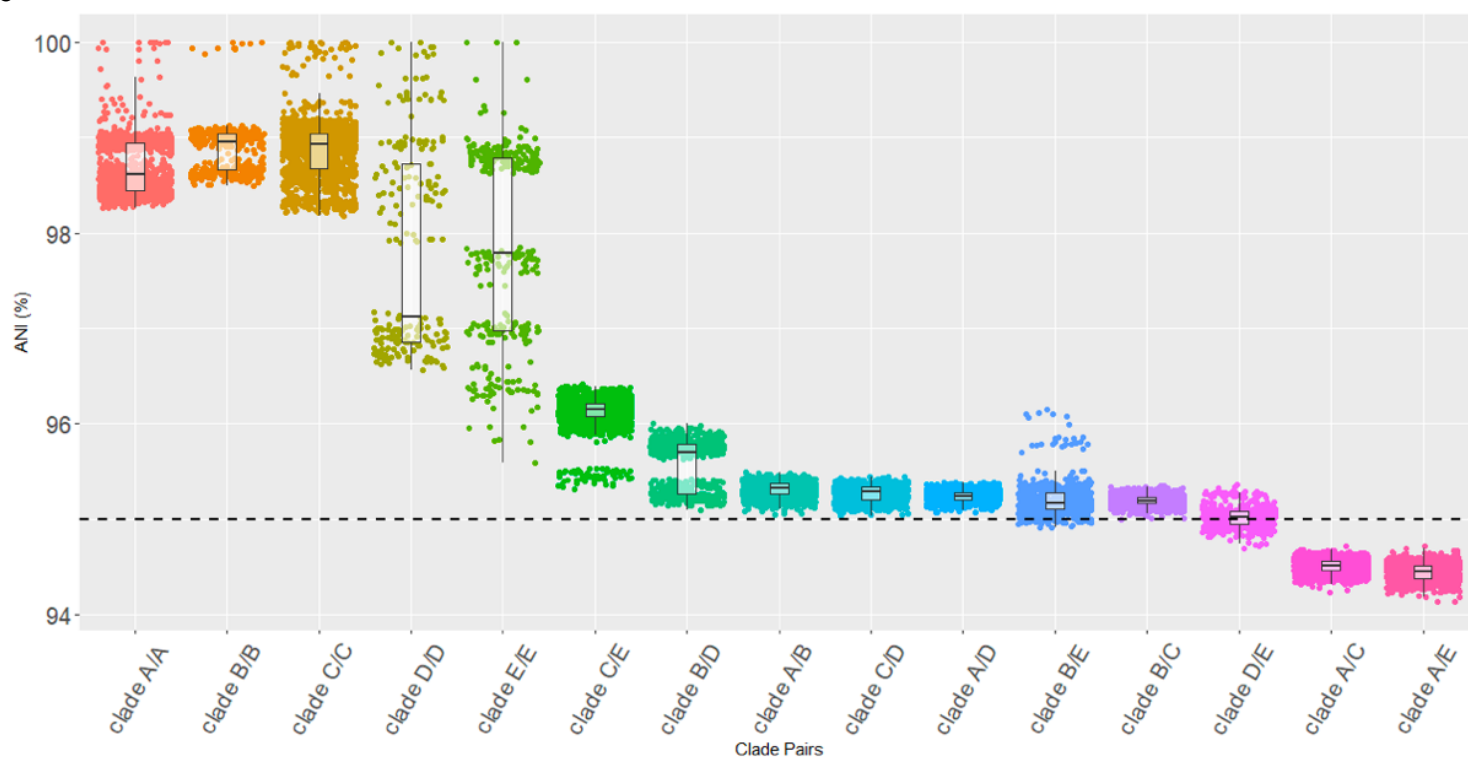
**Supplementary Fig. 2. Average nucleotide identity (ANI) comparisons of *P. aeruginosa* isolates.** Histogram of pairwise ANI values for 863 *P. aeruginosa* isolate genomes sampled by EDS-HAT. Dashed vertical line indicates 95% ANI. Comparisons in red are between isolates in *P. aeruginosa* Groups 1 or 2 versus isolates in the PA7-like Group 3, which appear to belong to a distinct genomospecies.
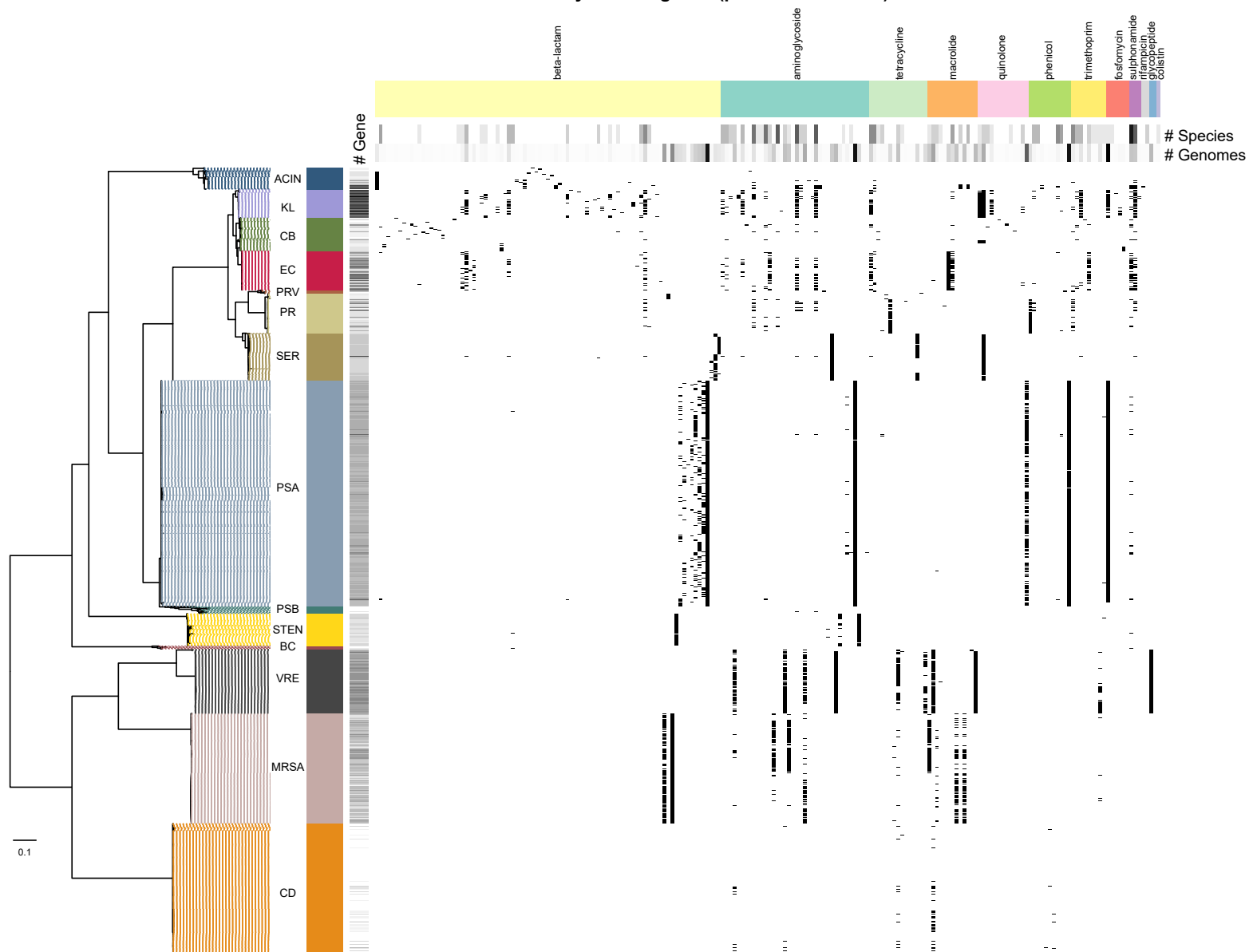
**Supplementary Fig. 3. Average nucleotide identity (ANI) comparisons of *S. marcescens* isolates. a**, Phylogeny and ANI of 177 *S. marcescens* isolates sampled by EDS-HAT. Grey shading indicates ANI values >95%, with darker shading showing higher identity. White indicates ANI values <95%. **b**, Distribution of pairwise ANI values for *S. marcescens* isolates belonging to the same or different clades, broken down into pairwise clade comparisons. All comparisons between isolates in Clade A vs. Clade C and Clade A vs. Clade E fall below the standard species cutoff of 95%.

**Supplementary Fig. 4. Distribution of antimicrobial resistance (AMR) genes among 3,004 clinical bacterial isolates from hospitalized patients.** Resistance genes were identified by BLASTn comparison to the ResFinder database. Isolates are ordered according to their phylogenetic placement using the amino acid sequences of 120 ubiquitous protein-coding genes from the Genome Taxonomy Database Tool Kit. "# Gene" shows the number of AMR genes per genome, with darker shading indicating more AMR genes. The matrix shows the presence or absence of 202 AMR genes, grouped by antibiotic class. Heat maps at the top show the number of species groups and total number of genomes encoding each gene, with darker shading indicating higher numbers. Raw data used to make the matrix are available in Table S3.