

1 **Empirical genomic methods for tracking plasmid spread among healthcare-associated**
2 **bacteria**

3

4 Running title: Epidemiology of plasmids in hospitals

5

6 Daniel Evans MS¹, Alexander Sundermann DrPH^{2,3}, Marissa Griffith BS^{2,3}, Vatsala Srinivasa

7 MPH^{2,3}, Mustapha Mustapha MBBS^{2,3}, Jieshi Chen MS⁴, Artur Dubrawski PhD⁴, Vaughn

8 Cooper PhD^{3,5,6}, Lee Harrison MD^{2,3}, Daria Van Tyne PhD^{2,3,6}

9

10 ¹Department of Infectious Diseases and Microbiology, University of Pittsburgh Graduate School
11 of Public Health, Pittsburgh PA, USA

12 ²Division of Infectious Diseases, University of Pittsburgh School of Medicine, Pittsburgh PA,
13 USA

14 ³Center for Genomic Epidemiology, University of Pittsburgh School of Medicine, Pittsburgh PA,
15 USA

16 ⁴Auton Laboratory, Carnegie Mellon University, Pittsburgh PA, USA

17 ⁵Department of Microbiology and Molecular Genetics, University of Pittsburgh School of
18 Medicine, Pittsburgh PA, USA

19 ⁶Center for Evolutionary Biology and Medicine, University of Pittsburgh School of Medicine,
20 Pittsburgh, PA, USA

21

22 Correspondence to Dr. Daria Van Tyne, PhD, Division of Infectious Diseases, University of

23 Pittsburgh School of Medicine, Pittsburgh PA 15213, USA; +1 412 648 4210; vantyne@pitt.edu

24 **Summary**

25 **Background:** Healthcare-associated bacterial pathogens frequently carry plasmids that
26 contribute to antibiotic resistance and virulence. The horizontal transfer of plasmids in healthcare
27 settings has been previously documented, but genomic and epidemiologic methods to study this
28 phenomenon remain underdeveloped. The objectives of this study were to develop a method to
29 systematically resolve and track plasmids circulating in a single hospital, and to identify
30 epidemiologic links that indicated likely horizontal plasmid transfer.

31 **Methods:** We derived empirical thresholds of plasmid sequence similarity from comparisons of
32 plasmids carried by bacterial isolates infecting individual patients over time, or involved in
33 hospital outbreaks. We then applied those metrics to perform a systematic screen of 3,074
34 genomes of nosocomial bacterial isolates from a single hospital for the presence of 89 plasmids.
35 We also collected and reviewed data from electronic health records for evidence of geotemporal
36 associations between patients infected with bacteria encoding plasmids of interest.

37 **Findings:** Our analyses determined that 95% of analyzed genomes maintained roughly 95% of
38 their plasmid genetic content at a nucleotide identity at least 99-985%. Applying these similarity
39 thresholds to identify horizontal plasmid transfer identified 45 plasmids circulating among
40 clinical isolates. Ten plasmids met criteria for geotemporal links associated with horizontal
41 transfer. Several plasmids with shared backbones also encoded different additional mobile
42 genetic element content, and these elements were variably present among the sampled clinical
43 isolate genomes.

44 **Interpretation:** The horizontal transfer of plasmids among nosocomial bacterial pathogens is
45 frequent within hospitals and can be monitored with whole genome sequencing and comparative

46 genomics approaches. These approaches should incorporate both nucleotide identity and
47 reference sequence coverage to study the dynamics of plasmid transfer in the hospital.
48 **Funding:** This research was supported by the US National Institute of Allergy and Infectious
49 Disease (NIAID) and the University of Pittsburgh School of Medicine.

50 **RESEARCH IN CONTEXT**

51 **Evidence before this study:**

52 A search of PubMed for research articles containing the search terms “plasmid”,
53 “transfer”, “epidemiology”, “hospital”, and “patients” identified 115 peer-reviewed manuscripts
54 published before 01 January 2022. Twenty-four manuscripts documented the dissemination of
55 one or more plasmids by horizontal transfer in a hospital setting. Most of these prior studies
56 focused on a single plasmid, outbreak, antibiotic resistance gene or pathogen species, and none
57 established an *a priori* approach to identify plasmids circulating among non-clonal bacterial
58 genomes. While prior studies have quantified plasmid preservation and nucleotide identity,
59 similarity thresholds to infer horizontal transfer were neither uniform across studies nor
60 systematically derived from empirical data.

61 **Added value of this study:**

62 This study advances the field of genomic epidemiology by proposing and demonstrating
63 the utility of empirically derived thresholds of plasmid sequence similarity for inferring
64 horizontal transfer in healthcare settings. It also advances the field by tracking horizontal plasmid
65 transfer within a single hospital at a hitherto unprecedented scale, examining the evidence of
66 horizontal transfer of 89 plasmids among thousands of clinical bacterial isolates sampled from a
67 single medical center. Our systematic review of patient healthcare data related to horizontal
68 transfer also occurred at a breadth not previously undertaken in hospital epidemiology.

69 **Implications of all the available evidence:**

70 When successfully integrated into contemporary methods for surveillance of nosocomial
71 pathogens, comparative genomics can be used to track and intervene directly against the
72 dissemination of plasmids that exacerbate virulence and antimicrobial resistance in healthcare-

73 associated bacterial infections. Standardized thresholds of plasmid identity benefit epidemiologic
74 investigations of horizontal transfer similar to those offered by establishing uniform thresholds of
75 genome identity for investigations of bacterial transmission.

76 **Introduction**

77 Healthcare-associated infections impose a serious burden on healthcare infrastructure and
78 widespread morbidity and mortality, both in the United States and worldwide.¹ Plasmids carried
79 by bacterial pathogens that cause these infections often carry genes conferring antimicrobial
80 resistance, virulence, and environmental persistence. These features complicate patient care and
81 increase disease severity.^{2,3} While our understanding of plasmid transmission via horizontal
82 transfer among bacteria in healthcare settings has increased in recent years, it remains poorly
83 understood compared to healthcare-associated transmission of bacteria.⁴⁻⁶ One substantial gap in
84 the burgeoning field of plasmid epidemiology is the lack of uniform thresholds of sequence
85 similarity by which recent horizontal plasmid transfer can be inferred. Without standardized
86 genomics-based approaches to resolve and characterize plasmids found in healthcare settings,
87 inferences regarding plasmid transmission remain speculative and error-prone. Developing
88 metrics to establish plasmid transmission can aid infection prevention personnel in the spread of
89 multidrug-resistant bacteria in hospitals, helping to intervene directly against the dissemination
90 of antimicrobial resistance and virulence.^{2,5,7}

91 Our study had two objectives. The first was to establish sequence similarity thresholds to
92 infer horizontal transfer of plasmids among bacteria collected from a single hospital. The second
93 was to apply these thresholds to a large dataset of whole-genome sequences of nosocomial
94 bacterial pathogens from a tertiary hospital,^{8,9} and to investigate epidemiologic links between
95 patients infected with genetically unrelated bacteria encoding the same plasmid. The intended
96 outcome was to demonstrate that these methods augment our knowledge of plasmid dynamics in
97 healthcare settings at the scale of an entire hospital, further supporting the value of whole-
98 genome sequencing (WGS) in healthcare epidemiology.⁹⁻¹¹

99

100 **Methods**

101 Collection of clinical bacterial isolates and corresponding patient data

102 The study took place at the University of Pittsburgh Medical Center (UPMC)
103 Presbyterian Hospital and the University of Pittsburgh School of Medicine. All bacterial isolates
104 analyzed in this study were collected through the Enhanced Detection System for Hospital-
105 Associated Transmission (EDS-HAT) project, using previously published eligibility criteria for
106 selection of isolates.^{8,9,11} EDS-HAT uses a combination of WGS surveillance of selected
107 hospital-associated pathogens and machine learning of the electronic health record to identify
108 otherwise undetected outbreaks and the responsible transmission routes, respectively. Isolates
109 were identified using TheraDoc software (Version 4.6, Premier, Inc, Charlotte NC). The
110 University of Pittsburgh Institutional Review Board (IRB) approved the EDS-HAT project and
111 classified it as being exempt from informed consent because direct contact with human subjects
112 is not performed. Approval was also granted (STUDY20060252) to use data collected for the
113 EDS-HAT project to study the horizontal transfer of mobile genetic elements in the hospital.

114

115 Short- and long-read whole-genome sequencing, processing, and assembly

116 DNA was extracted from overnight cultures of single bacterial colonies using a Qiagen
117 DNeasy Blood and Tissue Kit (Qiagen, Germantown MD). Short-read WGS was performed
118 using the Illumina platform (Illumina, San Diego CA), with libraries constructed with a Nextera
119 DNA Sample Prep Kit with 150bp paired-end reads and sequenced on the NextSeq platform.
120 Trim Galore v0.6.1 was used to remove sequencing adaptors, low-quality bases, and poor-quality
121 reads.¹² Bacterial species were identified from processed Illumina reads by alignment to Kraken
122 v1.0 and RefSeq databases.^{13,14} Genomes of strains sequenced only with Illumina technology

123 were assembled using SPAdes v3.11.¹⁵ The quality of assembled genomes was then verified
124 using QUAST.¹⁶ Assembled genomes were excluded if they failed to meet the following quality
125 control parameters: genome-wide read depth of at least 40X, cumulative length of assembled
126 genome within 20% of the expected length for the assigned genus, fewer than 400 contigs in the
127 assembled genome, and an N50 value of greater than 50,000bp. Illumina sequencing data for all
128 isolates is deposited at NCBI with accession numbers listed in **Appendix 1**.

129 Long-read sequencing and base-calling of select isolates were performed using Oxford
130 Nanopore technology (Oxford Nanopore Technologies, Oxford, United Kingdom). Libraries
131 were constructed using a rapid multiplex barcoding kit (catalog number SQK-RBK004).
132 Sequencing was performed using an Oxford Nanopore MinION device with R9.4.1 flow cells.
133 Base-calling and read processing was performed using Albacore v2.3.3 or Guppy v2.3.1 (Oxford
134 Nanopore Technologies, Oxford, United Kingdom) using default parameters. Hybrid assembly
135 was performed for genomes for which both short- and long-read sequencing data were available
136 and whose short-read only assemblies passed the aforementioned parameters, using Unicycler
137 v0.4.7 or v0.4.8-beta.¹⁷

138

139 Genome annotation, plasmid assembly, characterization, alignment, and phylogenetic analyses

140 Assembled genomes and hybrid-assembled plasmids were annotated with Prokka v1.13
141 or v1.14.¹⁸ Multi-locus sequence types (STs) were assigned using *mlst* v2.16.1 with PubMLST
142 typing schemes.^{19,20} Antibiotic resistance genes were identified by BLASTn alignments of
143 assembled genomes to the ResFinder v4.1 database.^{21,22} Plasmid replicons were identified by
144 BLASTn alignments of assembled genomes to the PlasmidFinder v2.1 database.²³ Other
145 genomic features were identified from annotations by Prokka¹⁸ and by BLASTn alignments to

146 the VFDB database.²⁴ Species were defined by grouping isolates with average nucleotide identity
147 (ANI) of at least 95% to one another and less than 95% ANI to genomes of other species.²⁵
148 Whole-genome phylogenies were constructed from core genome alignments generated by Roary
149 v5.18.2²⁶ using RAxML v8.0.26 with 1,000 bootstrap iterations.²⁷ Plasmid phylogenies were
150 constructed from core genome alignments of reference plasmids made with *snippy-core* v4.4.5²⁸
151 using RAxML v8.0.26 with 1,000 bootstrap iterations.²⁷ Annotated plasmids were aligned to one
152 another using EasyFig v2.2.2.²⁸ Plasmids were named according to the genus- or species-based
153 identification code of the isolate from which they were first identified and the contig in the
154 hybrid assembly corresponding to the closed plasmid sequence, separated by an underscore.

155

156 Sequence similarity analyses to establish metrics of plasmid transmission

157 Plasmids were identified from hybrid-assembled genomes of clonal isolates implicated in
158 outbreaks within the same hospital.^{10,11} Contigs were classified as plasmids if they were
159 circularized during hybrid assembly,¹⁷ were 2-300kb in length, and possessed at least one
160 replicon or conjugative plasmid gene identified by PlasmidFinder, Prokka, or RAST.^{18,23,30} The
161 proportion of plasmid gene content preserved in bacterial genomes (“sequence coverage”) was
162 calculated by searching for plasmid sequences among the contigs of each assembled bacterial
163 genome using BLASTn,³⁰ using a nucleotide sequence identity threshold of 95%. Total plasmid
164 nucleotide sequence identity was quantified based on single nucleotide polymorphisms (SNPs)
165 identified by *snippy-core* v4.4.5.²⁸

166

167 Resolving plasmids circulating among nosocomial bacterial pathogens in a single hospital

168 Eighty-nine closed reference plasmid sequences that met aforementioned sequence
169 coverage criteria were selected from whole-genome bacterial sequences constructed using
170 previously published methods (**Appendix 2**).^{5,9,11} Plasmids were de-duplicated by aligning and
171 visualizing closed reference plasmid sequences of similar lengths and distribution among STs
172 and species using Mauve or EasyFig v2.2.2.^{29,31} 3,074 whole-genome sequences from clinical
173 bacterial isolates were then screened for the presence of 89 plasmids by calculating sequence
174 coverage by BLASTn-based contig mapping³² and sequence identity using *snippy-core*,²⁸ as
175 described above (**Supplementary Figure 1**). Plasmids were considered potentially involved in
176 vertical transmission and/or horizontal transfer if they were present at sufficient nucleotide
177 similarity and sequence coverage in at least two bacterial genomes of the same or different
178 bacterial species.

179

180 Identification of geotemporal associations where horizontal plasmid transfer may have occurred

181 Geotemporal associations between plasmid-carrying isolates were identified by
182 systematic review of patient electronic health record (EHR) data, using a previously published
183 machine learning algorithm based on case-control methodology followed by manual review of
184 the algorithm's findings.³³ Geotemporal links highlighted by the algorithm were manually
185 evaluated if (1) they included bacterial isolates of different STs, species, or genera; (2) at least
186 one patient was infected during, after, or shortly before the date of culture of their plasmid-
187 carrying isolate while in the location; and (3) at least one other patient later infected by another
188 isolate with the same plasmid was exposed to the same location within 90 days of the date(s) of
189 culture of their isolate(s). Admission records for roommates with overlapping stays were
190 identified as shared exposures regardless of the length of cohabitation. Manual investigations of

191 patient clusters of interest were performed by a board-certified infection preventionist (AJS),
192 with the investigator blinded to genomic data regarding associated pathogens or plasmids.

193

194 **Results**

195 Empirical thresholds of plasmid similarity to indicate potential horizontal transfer

196 While the epidemiology of plasmid transfer among hospital-associated bacteria has
197 previously been described,^{4,5,34,35} the field lacks consistent guidelines for determining when
198 horizontal plasmid transfer among infection-derived bacterial isolates can be inferred. We
199 hypothesized that a nucleotide sequence similarity-based threshold could be developed by first
200 quantifying the sequence similarity of plasmids that were carried by bacterial strains known to
201 have been transmitted among hospital inpatients. We first examined four plasmids carried by
202 three clusters of carbapenemase-producing *Klebsiella pneumoniae* and vancomycin-resistant
203 *Enterococcus faecium* that showed strong genomic and epidemiologic evidence of nosocomial
204 transmission within our hospital.^{10,11} We also examined plasmids present in pairs or groups of
205 clinical isolates of identical species and sequence types that were isolated from the same patients
206 at different dates during their hospital stays. These plasmids were presumed to be vertically
207 transmitted during bacterial cell division and long-term carriage within each patient. In total, we
208 performed 57 comparisons between 25 plasmids in 47 bacterial genomes; 15 of these
209 comparisons were between same-patient isolates. Elapsed time between dates of culture of
210 plasmid-linked strains ranged from zero days to 427 days (mean 119.5 days, median 82 days).

211 We next used a previously published BLASTn-based method to calculate the proportion
212 of each plasmid's sequence that was preserved in linked isolates.⁵ We also quantified the number
213 of single nucleotide polymorphisms (SNPs) among the shared plasmid sequences present in each
214 group of epidemiologically linked isolates. We refer to these metrics as "plasmid sequence

215 coverage” and “plasmid sequence identity”, respectively (**Figure 1a**). Ninety-five percent of
216 these comparisons had sequence coverage values of at least 93.7% (**Figure 1b**) and nucleotide
217 sequence identity of at least 99.985% (i.e. fewer than 15 SNPs per 100kb of plasmid sequence,
218 or 1 SNP per 6.67kb) (**Figure 1c**). These results provided initial estimates of sequence similarity
219 thresholds that were then applied for downstream analyses to detect horizontal transfer within the
220 hospital environment.

221

222 Systematic identification of plasmid horizontal transfer from a single hospital

223 To systematically identify potential horizontal transfer of plasmids within our hospital,
224 we applied our empirical thresholds described above to search for the presence of 89 reference
225 plasmids in the genomes of 3,074 clinical bacterial isolates collected from 2,086 patients
226 (**Supplementary Figure 1**) (**Appendix 2**). Twelve of these plasmids were previously described
227 in a prior study of mobile genetic element diversity in our hospital.⁵ We identified and additional
228 34 plasmids carried by isolates from more than one patient that had sequence coverage $\geq 95\%$
229 and sequence identity $\geq 99.985\%$ compared to the reference plasmid sequence. The 46 shared
230 plasmids were present in 336 of 3,074 (10.9%) bacterial isolates and in seven of 12 (58.3%)
231 genera in our genomic dataset (**Table 1**). Bacterial isolates carrying shared plasmids were
232 collected from 291 of 2,086 (14%) unique patients. Shared plasmids were most abundant among
233 isolates belonging to *Enterococcus* (19 plasmids in 214 isolates, 85.6% of *Enterococcus* isolates
234 in the dataset) (**Figure 2a**) and *Klebsiella* (16 plasmids in 77 isolates, 56.6%) (**Figure 2b**). Nine
235 of the 46 shared plasmids were detected in isolates belonging to different species, and 23 others
236 were found in isolates of different multi-locus STs within the same species, suggesting horizontal
237 plasmid transfer (**Appendix 2**). The remaining 15 plasmids were found in isolates belonging to

238 the same species and ST, suggesting they may have been vertically transmitted between patients
239 along with the bacteria carrying them. Taken together, these results demonstrate the widespread
240 distribution of numerous plasmids engaged in vertical transmission and horizontal transfer within
241 a single hospital.

242

243 Geotemporal associations among patients implicated in horizontal plasmid transfer

244 To investigate the epidemiology of horizontal plasmid transfer, we applied a previously
245 published screening algorithm³³ to systematically review hospital electronic health record (EHR)
246 data of patients infected with isolates belonging to different STs, species, or genera that shared
247 plasmids with high nucleotide identity and high sequence coverage. Using this approach, we
248 identified 18 groups of patients with geotemporal associations who were infected with isolates
249 encoding eight different plasmids that were likely horizontally transferred (**Table 2, Appendix**
250 **2**). These associations fell into three major categories, in which previously uninfected patients
251 were found to be culture-positive for plasmid-carrying pathogens after they were: (1) admitted to
252 rooms previously occupied by infected patients, (2) admitted to units simultaneously occupied by
253 at least one infected patient, or (3) admitted to units previously occupied by at least one infected
254 patient. These findings suggest that plasmids are frequently horizontally transferred among
255 pathogens infecting hospitalized patients independent from bacterial transmission.

256

257 Nucleotide sequence identity threshold improves identification of horizontally transferred 258 plasmids

259 To investigate whether the methods we developed provided improved resolution over our
260 previously published sequence coverage-based approach,⁵ we reanalyzed ten plasmids that were

261 resolved and investigated in our earlier study. We found that a small ColRNAI plasmid
262 (pKLP00155_6, 9.5kb) that was detected in a large number of isolates of multiple species was
263 most likely two closely related but distinct plasmids that were co-circulating in our hospital
264 (**Figure 3a**). Pairwise comparisons to the pKLP00155_6 reference sequence revealed a group of
265 13 isolates with >95% sequence coverage of the reference plasmid, but whose sequence
266 identities ranged from 42 to 73 SNPs per 100kb. When the same sequences were compared to
267 pCB00073_2, another reference plasmid that had previously been considered the same as
268 pKLP00155_6 based on sequence coverage alone,⁵ their pairwise sequence identities ranged
269 from 0 to 11 SNPs per 100kb. Separation of these two plasmids was further validated by
270 phylogenetic analysis, which clearly delineated the plasmid sequences into two groups (**Figure**
271 **3b**). When we investigated geotemporal associations between patients infected with bacteria
272 carrying these two plasmids (**Table 2**), we found that two patients infected by isolates of
273 different bacterial species carrying pCB00073_2 were both cultured during overlapping stays on
274 the same hospital unit (**Figure 3c**). Unit-based horizontal transmission was also observed for 16
275 of the 24 patients carrying pKLP00155_6 housed in three different units, with one patient linked
276 to two of these units (**Figure 3d**).

277 We also observed several cases of plasmid pairs that met the sequence identity threshold
278 of $\geq 99.985\%$ identity, but did not meet the $\geq 95\%$ sequence coverage threshold. Sequence
279 alignments confirmed that these plasmid pairs each shared a core or “backbone” plasmid
280 sequence, with the longer plasmid of each pair carrying additional genetic cargo of transposases
281 and/or IS6, IS26, and IS110-mediated insertion sequences (**Supplementary Figure 2**).

282 Phylogenetic analysis of one of these pairs, pKLP00161_2 and pKLP00218_2, did not separate
283 the plasmids into distinct clades defined by the presence or absence of the additional MGE,
284 suggesting potential gain or loss of plasmid cargo during co-temporal circulation of the bacterial

285 isolates carrying these plasmids (**Supplementary Figure 3**). We considered isolates with high
286 sequence identity and high sequence coverage to one or the other of these plasmid pairs as
287 encoding the same plasmid backbone for the purposes of assessing epidemiologic links between
288 patients.

289 We next performed extensive reviews of 19 patient care records associated with
290 pathogens of different STs or species that carried four different plasmid variants with shared
291 backbones. Ten patients infected with multiple STs of *K. pneumoniae* and one ST of *K. oxytoca*
292 that all carried the pKLP00218_2 backbone sequence had potential exposures across four
293 different hospital units (**Figure 4a**). Nine patients infected with five different STs of *E. faecium*
294 carrying the pVRE32553_4 backbone sequence had potential exposures across three different
295 hospital units (**Figure 4b**). Taken together, these findings demonstrate how empirically derived
296 plasmid similarity thresholds enable tracking of plasmid circulation and horizontal transfer
297 among bacteria. Moreover, these data highlight the potential plasticity of plasmid backbones as
298 they circulate in a hospital.

299

300 **DISCUSSION**

301 In this study, we developed uniform thresholds of genetic similarity among plasmids
302 from which horizontal transfer in a hospital might be inferred. We then applied these thresholds
303 to a large dataset of whole-genome sequences of nosocomial bacterial isolates collected from
304 patients at a single medical center. We observed that plasmids potentially engaged in horizontal
305 transfer were widespread among infected patients, and subsequent reviews of electronic health
306 record data found numerous geotemporal links that may have provided opportunities for plasmid
307 transfer between patients. Our findings extend our previous exploration into the genomics of

308 horizontal plasmid transfer within hospital settings,⁵ in which we identified likely plasmid
309 transfer among bacterial pathogens that either co-infected patients or carried those plasmids
310 during transmission between patients with shared geotemporal risk factors.

311 This study contributes to the field of bacterial genomic epidemiology in several key
312 ways. While numerous methods have been previously applied to resolve and characterize
313 plasmids present among environmental and clinical bacterial isolates from hospitals,^{4,5} these
314 prior studies have largely focused on transmission of antibiotic resistance elements^{4,36} or
315 individual plasmids,^{34,35} often as part of corollary findings in studies focused on bacterial
316 transmission.⁶ In contrast, here we developed and implemented an empirical genome sequence-
317 based approach to systematically cluster plasmids from nosocomial isolates *en masse*. This study
318 also lays the foundation for more detailed investigations into the mechanics of horizontal
319 plasmid transfer within clinical settings^{2,35}. Lastly, our work further supports the utility of
320 comprehensive whole-genome sequencing of bacterial isolates from the hospital.^{4,5,34}

321 Our findings bolster several recent insights on the potential causes and effects of plasmid
322 sharing among bacteria infecting hospitalized patients. For example, only 45 of the 89 reference
323 plasmids we studied carried known genes conferring antimicrobial resistance (**Appendix 2**). This
324 finding is in agreement with recent research suggesting that while plasmids mediate much of the
325 spread of antimicrobial resistance, antimicrobial resistance itself may not be the primary driver
326 of horizontal plasmid transfer in clinical settings.³⁷ This study also underscores the need for a
327 more thorough characterization of the plasmidome, as the majority of genes on sequenced
328 plasmids have unknown functions.³⁸ We also observed widespread sharing of several plasmids
329 among vancomycin-resistant *Enterococcus* (VRE) isolates, corroborating previous findings on
330 the impact of the diversity and extent of plasmid dissemination on enterococcal evolution.³⁹

331 Even with the increased scope of our genomic and epidemiologic analysis, we postulate
332 that our findings still underestimate the true extent to which horizontal plasmid transfer occurs
333 among nosocomial bacterial pathogens. The large number of hybrid-assembled reference
334 plasmids we used to search for evidence of plasmid sharing were sequenced from only 56
335 isolates from a single hospital,⁸ representing a small subset of the bacteria that infected
336 hospitalized patients during the two-year period of this study. Within this relatively limited
337 dataset, we nonetheless identified several plasmids that exhibited evidence of further
338 recombination and rearrangement while circulating among patients in the hospital. Future studies
339 should further refine the methods we have developed here through analysis of additional
340 reference plasmids as well as isolates sampled from environmental sites or patient colonization.
341 Additionally, computational approaches that resolve plasmids from whole-genome sequence data
342 through non-reference-based approaches^{40,41} could identify additional cases of likely horizontal
343 plasmid transfer.

344 Our study had several limitations. Our sequence similarity criteria were developed by
345 analyzing a relatively small number of isolates and reference plasmids from a single hospital,
346 which may have introduced selection bias when determining numerical thresholds. Additionally,
347 our dataset of bacterial genomes was subject to selection criteria that limited the sampling of
348 some clinical isolates based on phenotypic antimicrobial resistance profiles and/or body sites.^{9,33}
349 The reference plasmids we studied were largely drawn from sequencing data of select pathogens
350 generated in previous projects,^{5,9,42} rather than from comprehensive efforts that may have yielded
351 more representative samples of plasmids within our hospital. To improve the specificity of our
352 epidemiologic analyses, we limited our scope to geotemporal associations among isolates with
353 substantial taxonomic distance. Lastly, our criteria for inferring horizontal transfer were based on

354 professional expertise in hospital infection prevention and the dynamics of bacterial transmission
355 among hospitalized patients, as knowledge of *in vivo* patterns of plasmid transfer in hospitals is
356 highly limited.^{5,35}

357 In summary, our study advances the field of bacterial genomic epidemiology by applying
358 empirically derived similarity thresholds encompassing both genetic content and nucleotide
359 identity to study plasmid spread among bacterial pathogens. The application of our methods to a
360 large genomic dataset from a single hospital resolved horizontal plasmid transfer at an
361 unprecedented scale and demonstrated the added value of routine whole-genome sequencing of
362 nosocomial bacterial pathogens in hospital infection prevention.^{8,9,35} Applying these refined,
363 plasmid-focused approaches in healthcare settings could help infection prevention teams better
364 understand signs and risk factors for the dissemination of drug resistance, allowing them to
365 improve their approaches to reduce the risk of severe infections among patients.^{5,7,35} Future work
366 in this area should focus on refining similarity thresholds and surveillance methods, further
367 examining patterns of mutation and genetic plasticity among circulating plasmids, and
368 implementing these methods in real-time to control the spread of plasmid-borne antimicrobial
369 resistance and virulence factors that exacerbate nosocomial infections.

370

371 **DATA SHARING**

372 Genome sequence data of all bacterial isolates studied here has been deposited in the
373 Sequence Read Archive (SRA) with accession numbers listed in Appendix 1. Plasmid sequences
374 included in this study have been deposited in GenBank with accession numbers listed in
375 Appendix 2.

376

377 **DECLARATION OF INTERESTS**

378 DE received compensation from the Allegheny County Health Department (Pittsburgh,
379 PA, USA) for services rendered as an infectious disease epidemiologist while performing
380 research and writing of this manuscript. DE, AS, and MM received compensation from
381 EpiCenter Genomics LLC (Pittsburgh, PA, USA) for consulting services in genomic
382 epidemiology. These entities did not provide funding directly towards the production of this
383 manuscript, nor did they pay these authors for work contributing to its publication.

384

385 **ACKNOWLEDGEMENTS**

386 We thank Jane Marsh, Daniel Snyder, Chinelo Ezeonwuka, and Kady Waggle for
387 assistance with clinical isolate collection and whole-genome sequencing. We also thank Melissa
388 Saul for curating electronic health record data for epidemiologic investigations, as well as Alecia
389 Rokes for assistance in curating and classifying reference plasmid sequences. Lastly, we thank
390 Kyle Miller for support in algorithm development and implementation.

391 REFERENCES

- 392 1 Centers for Disease Control and Prevention (CDC). 2019 National and State Healthcare-
393 Associated Infections (HAI) Progress Report. 2019. [https://arpsp.cdc.gov/profile/national-](https://arpsp.cdc.gov/profile/national-progress/united-states)
394 [progress/united-states](https://arpsp.cdc.gov/profile/national-progress/united-states).
- 395 2 Lerminiaux NA, Cameron ADS. Horizontal transfer of antibiotic resistance genes in clinical
396 environments. *Can J Microbiol* 2019; **65**: 34–44.
- 397 3 San Millan A. Evolution of Plasmid-Mediated Antibiotic Resistance in the Clinical Context.
398 *Trends Microbiol* 2018; **26**: 978–85.
- 399 4 Peter S, Bosio M, Gross C, *et al*. Tracking of Antibiotic Resistance Transfer and Rapid
400 Plasmid Evolution in a Hospital Setting by Nanopore Sequencing. *mSphere* 2020; **5**: e00525-
401 20, [/msphere/5/4/mSphere525-20.atom](https://doi.org/10.1128/mSphere.00525-20).
- 402 5 Evans DR, Griffith MP, Sundermann AJ, *et al*. Systematic detection of horizontal gene
403 transfer across genera among multidrug-resistant bacteria in a single hospital. *eLife* 2020; **9**.
404 DOI:10.7554/eLife.53886.
- 405 6 Tofteland S, Naseer U, Lislevand JH, Sundsfjord A, Samuelsen O. A long-term low-frequency
406 hospital outbreak of KPC-producing *Klebsiella pneumoniae* involving Intergenous plasmid
407 diffusion and a persisting environmental reservoir. *PloS One* 2013; **8**: e59015.
- 408 7 Salamzade R, Manson AL, Walker BJ, *et al*. Inter-species geographic signatures for tracing
409 horizontal gene transfer and long-term persistence of carbapenem resistance. *Genome Med*
410 2022; **14**: 37.
- 411 8 Mustapha MM, Srinivasa VR, Griffith MP, *et al*. Genomic Diversity of Hospital-Acquired
412 Infections Revealed through Prospective Whole-Genome Sequencing-Based Surveillance.
413 *mSystems* 2022; **7**: e01384-21.
- 414 9 Sundermann AJ, Chen J, Kumar P, *et al*. Whole-Genome Sequencing Surveillance and
415 Machine Learning of the Electronic Health Record for Enhanced Healthcare Outbreak
416 Detection. *Clin Infect Dis* 2021; : ciab946.
- 417 10 Sundermann AJ, Babiker A, Marsh JW, *et al*. Outbreak of Vancomycin-resistant
418 *Enterococcus faecium* in Interventional Radiology: Detection Through Whole Genome
419 Sequencing-Based Surveillance. *Clin Infect Dis Off Publ Infect Dis Soc Am* 2019; published
420 online July 16. DOI:10.1093/cid/ciz666.
- 421 11 Marsh JW, Mustapha MM, Griffith MP, *et al*. Evolution of Outbreak-Causing Carbapenem-
422 Resistant *Klebsiella pneumoniae* ST258 at a Tertiary Care Hospital over 8 Years. *mBio* 2019;
423 **10**. DOI:10.1128/mBio.01945-19.
- 424 12 Krueger F, James F, Ewels P, Afyounian E, Schuster-Boeckler B. FelixKrueger/TrimGalore:
425 v0.6.7 - DOI via Zenodo. 2021; published online July 23. DOI:10.5281/ZENODO.5127899.
- 426 13 Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
427 alignments. *Genome Biol* 2014; **15**: R46.
- 428 14 Pruitt KD, Tatusova T, Maglott DR. NCBI reference sequences (RefSeq): a curated non-
429 redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res*
430 2007; **35**: D61-65.
- 431 15 Bankevich A, Nurk S, Antipov D, *et al*. SPAdes: a new genome assembly algorithm and its
432 applications to single-cell sequencing. *J Comput Biol J Comput Mol Cell Biol* 2012; **19**: 455–
433 77.
- 434 16 Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome
435 assemblies. *Bioinforma Oxf Engl* 2013; **29**: 1072–5.
- 436 17 Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies
437 from short and long sequencing reads. *PLoS Comput Biol* 2017; **13**: e1005595.
- 438 18 Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinforma Oxf Engl* 2014; **30**:
439 2068–9.
- 440 19 Jolley KA, Bray JE, Maiden MCJ. Open-access bacterial population genomics: BIGSdb

- 441 software, the PubMLST.org website and their applications. *Wellcome Open Res* 2018; **3**:
442 124.
- 443 20 Seemann T. mlst. <https://github.com/tseemann/mlst>.
- 444 21 Bortolaia V, Kaas RS, Ruppe E, *et al.* ResFinder 4.0 for predictions of phenotypes from
445 genotypes. *J Antimicrob Chemother* 2020; **75**: 3491–500.
- 446 22 Zankari E, Hasman H, Cosentino S, *et al.* Identification of acquired antimicrobial resistance
447 genes. *J Antimicrob Chemother* 2012; **67**: 2640–4.
- 448 23 Carattoli A, Zankari E, García-Fernández A, *et al.* In silico detection and typing of plasmids
449 using PlasmidFinder and plasmid multilocus sequence typing. *Antimicrob Agents Chemother*
450 2014; **58**: 3895–903.
- 451 24 Liu B, Zheng D, Jin Q, Chen L, Yang J. VFDB 2019: a comparative pathogenomic platform
452 with an interactive web interface. *Nucleic Acids Res* 2019; **47**: D687–92.
- 453 25 Kim M, Oh H-S, Park S-C, Chun J. Towards a taxonomic coherence between average
454 nucleotide identity and 16S rRNA gene sequence similarity for species demarcation of
455 prokaryotes. *Int J Syst Evol Microbiol* 2014; **64**: 346–51.
- 456 26 Page AJ, Cummins CA, Hunt M, *et al.* Roary: rapid large-scale prokaryote pan genome
457 analysis. *Bioinformatics* 2015; **31**: 3691–3.
- 458 27 Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
459 phylogenies. *Bioinforma Oxf Engl* 2014; **30**: 1312–3.
- 460 28 Seemann, Torsten T. Snippy: rapid haploid variant calling and core SNP phylogeny.
461 Available. 2015. <https://github.com/tseemann/snippy>.
- 462 29 Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer. *Bioinforma*
463 *Oxf Engl* 2011; **27**: 1009–10.
- 464 30 Aziz RK, Bartels D, Best AA, *et al.* The RAST Server: Rapid Annotations using Subsystems
465 Technology. *BMC Genomics* 2008; **9**: 75.
- 466 31 Darling ACE. Mauve: Multiple Alignment of Conserved Genomic Sequence With
467 Rearrangements. *Genome Res* 2004; **14**: 1394–403.
- 468 32 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*
469 *Mol Biol* 1990; **215**: 403–10.
- 470 33 Sundermann AJ, Miller JK, Marsh JW, *et al.* Automated data mining of the electronic health
471 record for investigation of healthcare-associated outbreaks. *Infect Control Hosp Epidemiol*
472 2019; **40**: 314–9.
- 473 34 Prussing C, Snavely EA, Singh N, *et al.* Nanopore MinION Sequencing Reveals Possible
474 Transfer of blaKPC–2 Plasmid Across Bacterial Species in Two Healthcare Facilities. *Front*
475 *Microbiol* 2020; **11**: 2007.
- 476 35 R-GNOSIS WP5 Study Group, León-Sampedro R, DelaFuente J, *et al.* Pervasive
477 transmission of a carbapenem resistance plasmid in the gut microbiota of hospitalized
478 patients. *Nat Microbiol* 2021; **6**: 606–16.
- 479 36 Cerqueira GC, Earl AM, Ernst CM, *et al.* Multi-institute analysis of carbapenem resistance
480 reveals remarkable diversity, unexplained mechanisms, and limited clonal outbreaks. *Proc*
481 *Natl Acad Sci U S A* 2017; **114**: 1135–40.
- 482 37 Lopatkin AJ, Huang S, Smith RP, *et al.* Antibiotics as a selective driver for conjugation
483 dynamics. *Nat Microbiol* 2016; **1**: 16044.
- 484 38 Acman M, van Dorp L, Santini JM, Balloux F. Large-scale network analysis captures
485 biological features of bacterial plasmids. *Nat Commun* 2020; **11**: 2452.
- 486 39 Arredondo-Alonso S, Top J, McNally A, *et al.* Plasmids Shaped the Recent Emergence of the
487 Major Nosocomial Pathogen *Enterococcus faecium*. *mBio* 2020; **11**: e03284-19,
488 /mbio/11/1/mBio.03284-19.atom.
- 489 40 Laczny CC, Galata V, Plum A, Posch AE, Keller A. Assessing the heterogeneity of in silico
490 plasmid predictions based on whole-genome-sequenced clinical isolates. *Brief Bioinform*
491 2019; **20**: 857–65.

- 492 41Robertson J, Nash JHE. MOB-suite: software tools for clustering, reconstruction and typing of
493 plasmids from draft assemblies. *Microb Genomics* 2018; **4**. DOI:10.1099/mgen.0.000206.
494 42Babiker A, Evans DR, Griffith MP, *et al*. Clinical and Genomic Epidemiology of Carbapenem-
495 Nonsusceptible *Citrobacter* spp. at a Tertiary Health Care Center over 2 Decades. *J Clin*
496 *Microbiol* 2020; **58**: e00275-20.

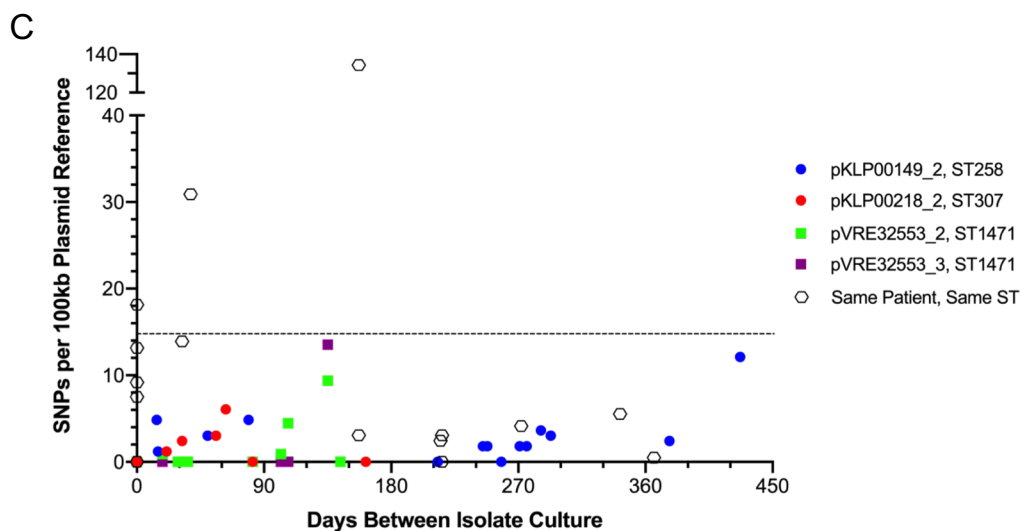
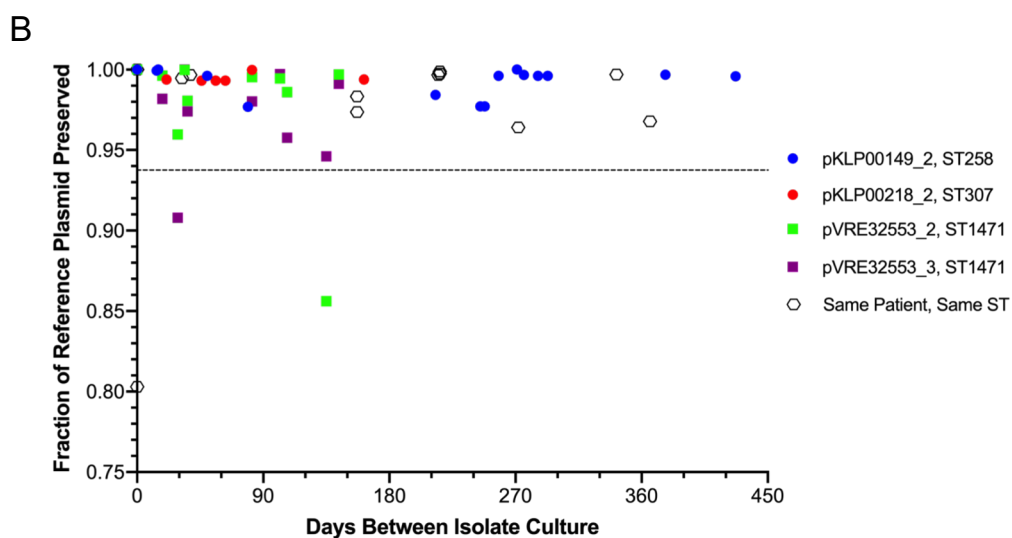
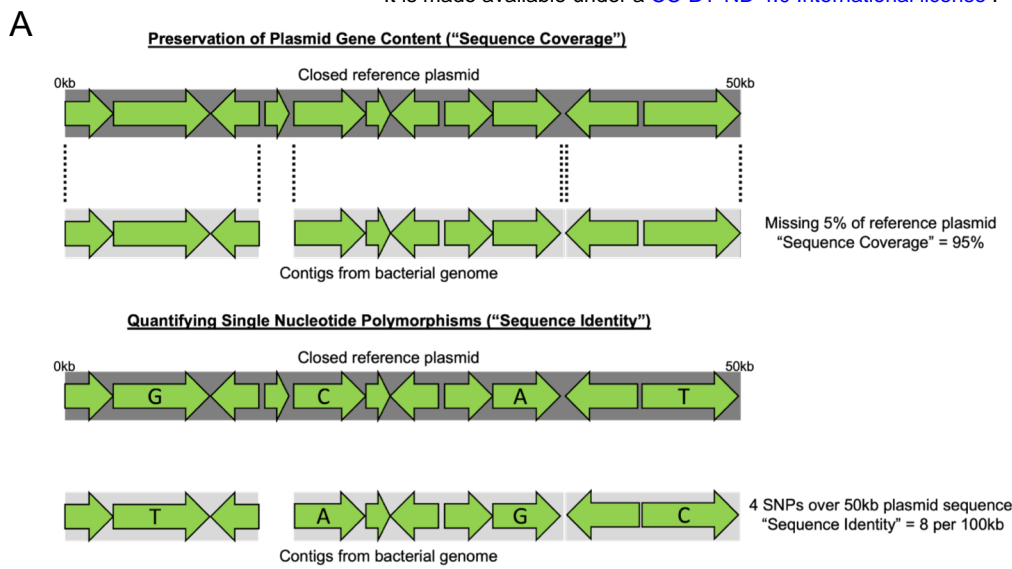
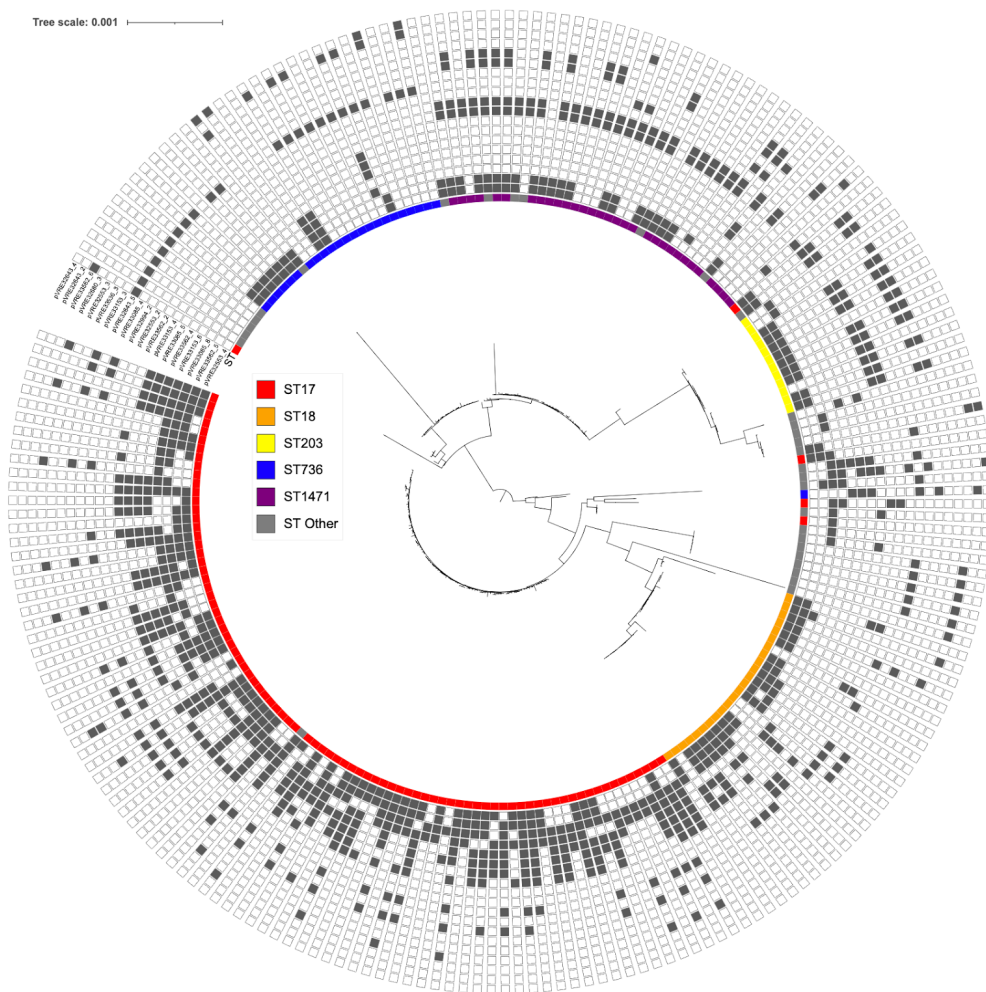


Figure 1: Empirically derived metrics of plasmid sequence similarity to infer horizontal plasmid transfer

(A) Approaches used to calculate (top) plasmid sequence coverage and (bottom) plasmid sequence identity. (B) Sequence coverage of plasmids carried by bacteria engaged in transmission between patients or prolonged carriage in the same patient. Dotted line shows the 5th percentile of coverage among analyzed plasmids. (C) Sequence identity of plasmids carried by isolates engaged in transmission between patients or prolonged carriage in the same patient. Dotted line shows the 95th percentile of SNPs per 100kb of plasmid reference sequence.

A



B

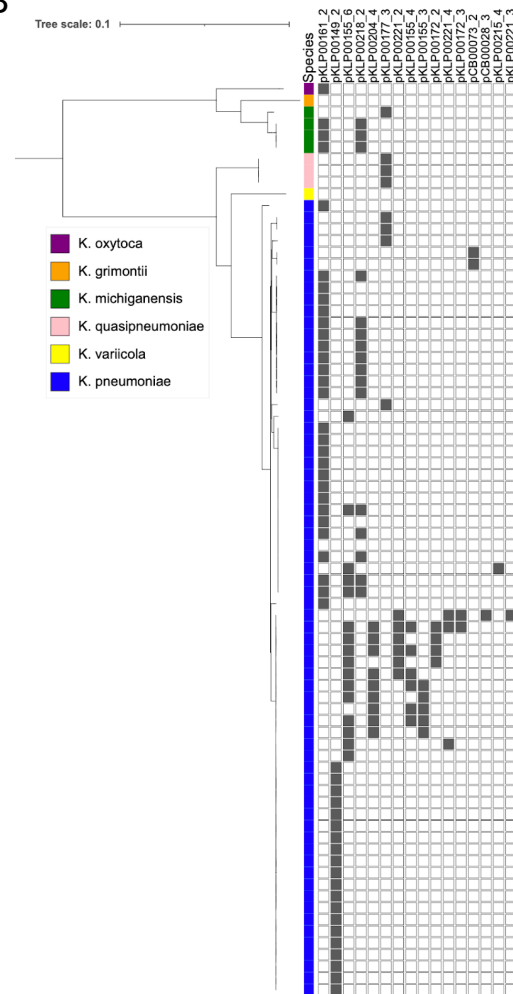


Figure 2: Plasmid distribution among nosocomial pathogens

Bootstrapped RAxML phylogenies of (A) vancomycin-resistant *Enterococcus faecium* and (B) *Klebsiella* spp. carrying shared plasmids. Color strips denote (A) sequence type (ST) of *E. faecium* or (B) species of *Klebsiella*.

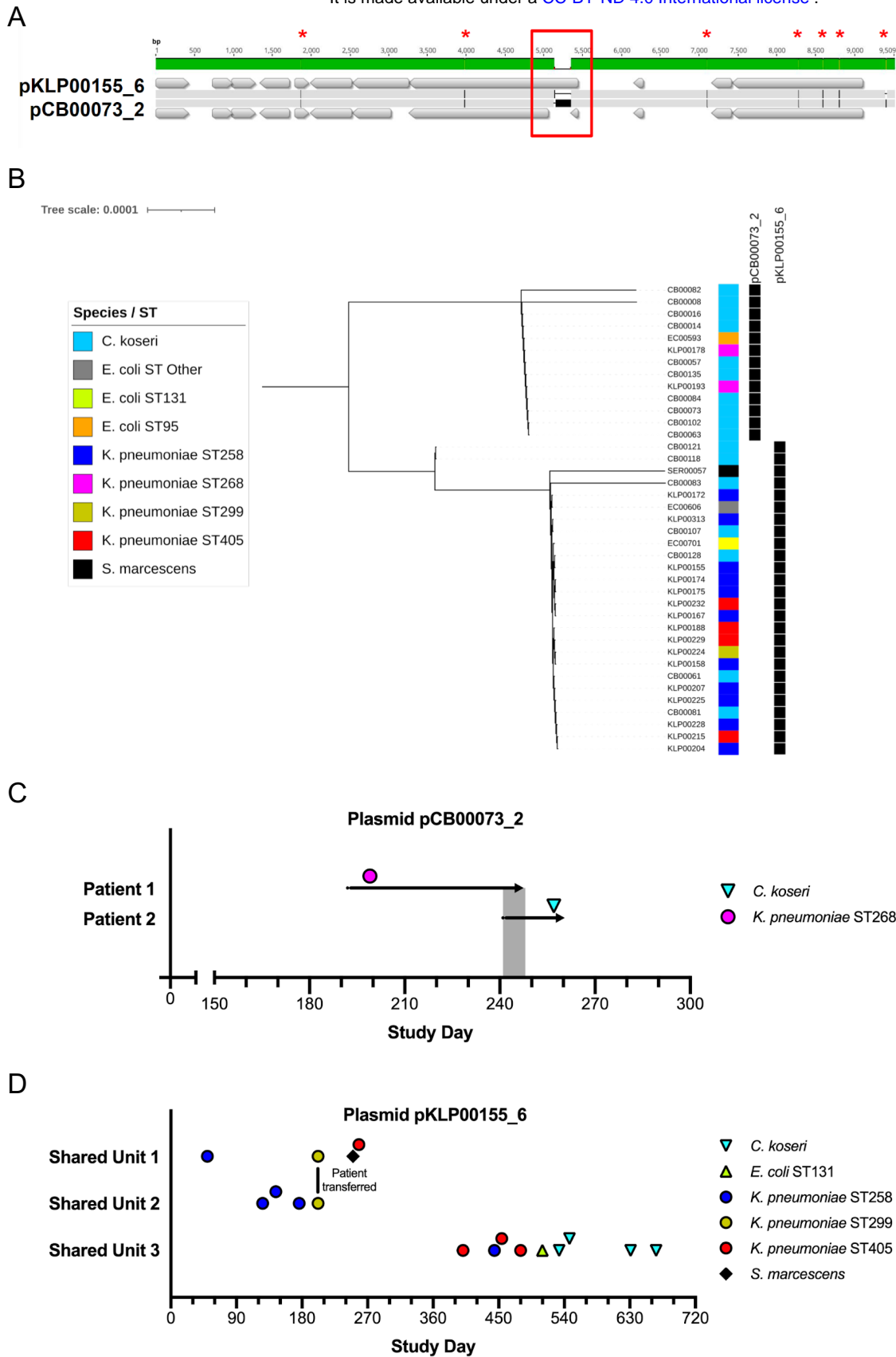
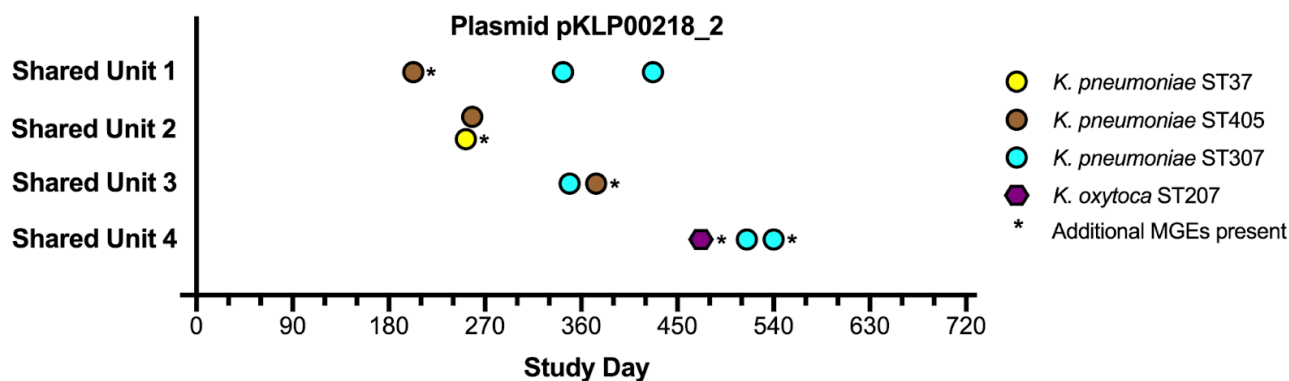


Figure 3: Resolution of two closely related colicin-encoding plasmid variants

(A) Nucleotide sequence alignment with point mutations highlighted by red asterisks and an insertion sequence marked by a red box. (B) Bootstrapped RAxML phylogeny of core plasmid sequences, using the pCB00073_2 sequence as an internal reference. (C) Timeline of dates of culture (de-identified as study days) and inpatient stays (gray shading) of two patients on the same unit infected with bacterial isolates of different species carrying pCB00073_2. (D) Timeline of culture dates of three groups of patients on three different shared units whose nosocomial isolates all carried pKLP00155_6.

A



B

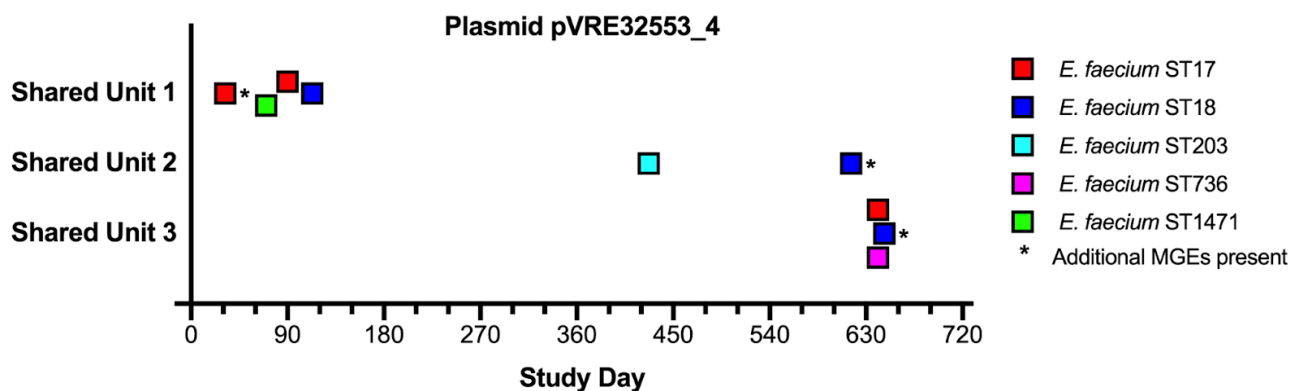
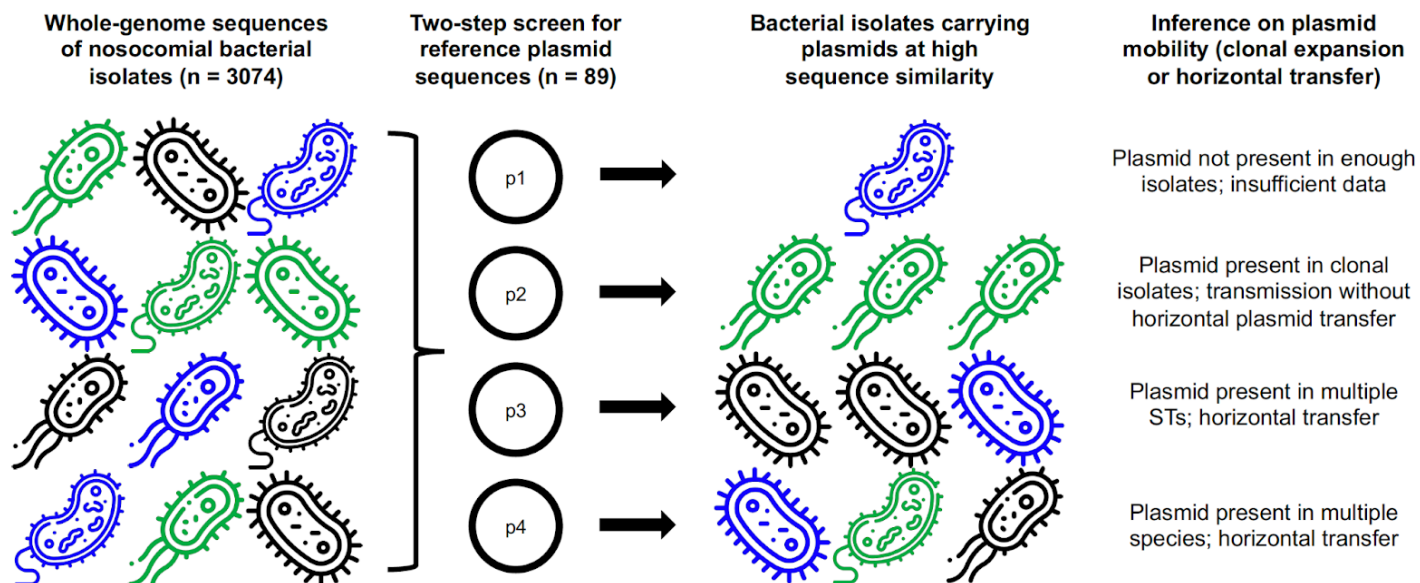
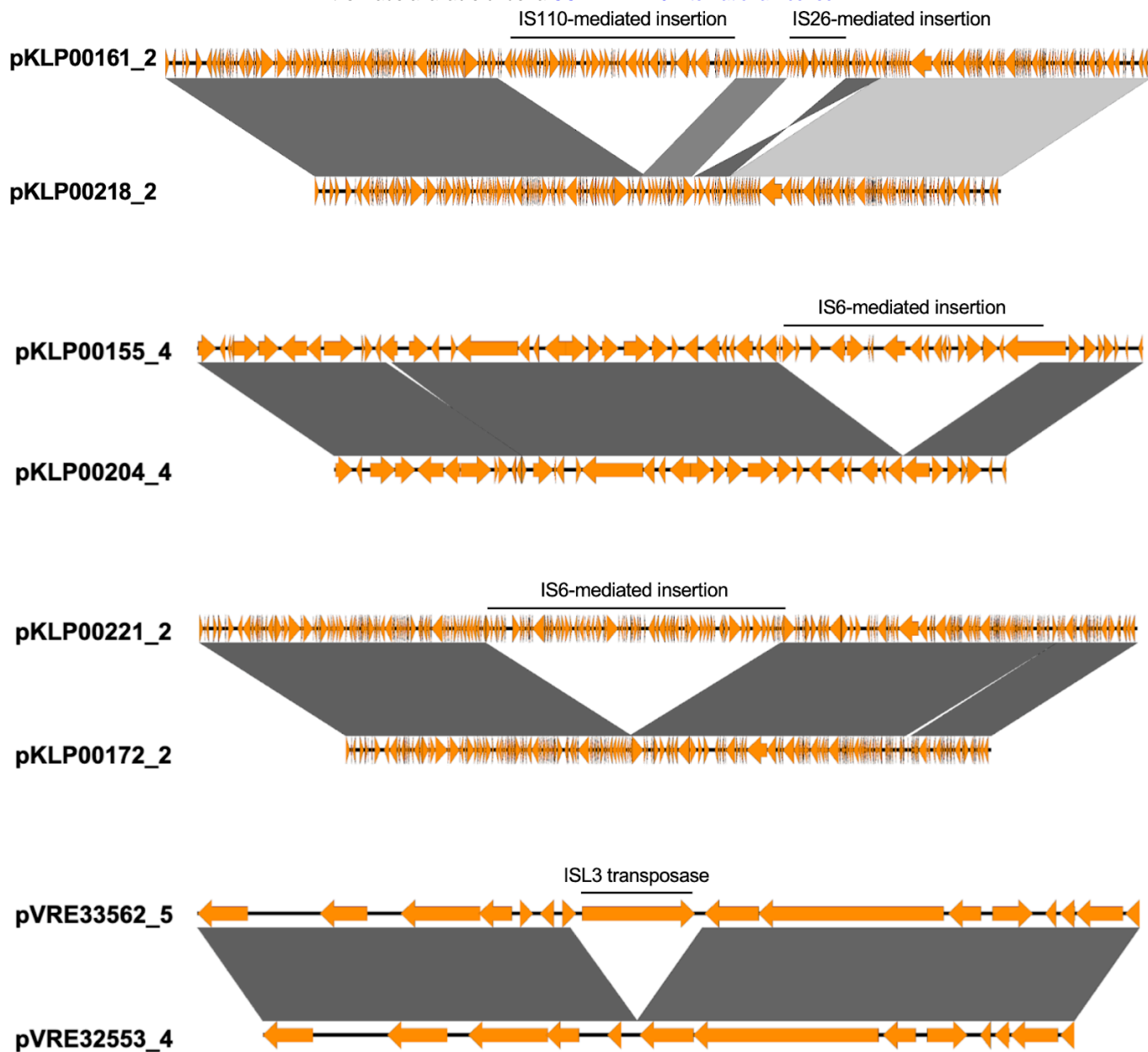


Figure 4: Geographically associated patients infected by bacteria carrying plasmid backbones engaged in horizontal transfer
Plots show associations of patients infected by (A) multiple STs and species of *Klebsiella* spp. carrying plasmid pKLP00218_2, and (B) multiple STs of *E. faecium* carrying pVRE32553_4. Study days correspond to dates of culture of plasmid-carrying pathogens. Shapes and colors of data points correspond to bacterial species and ST, respectively. See Supplementary Figure 2 for plasmid alignments depicting additional mobile genetic element (MGE) sequences.



Supplementary Figure 1: Conceptual flow diagram of alignment-based pipeline to infer horizontal plasmid transfer

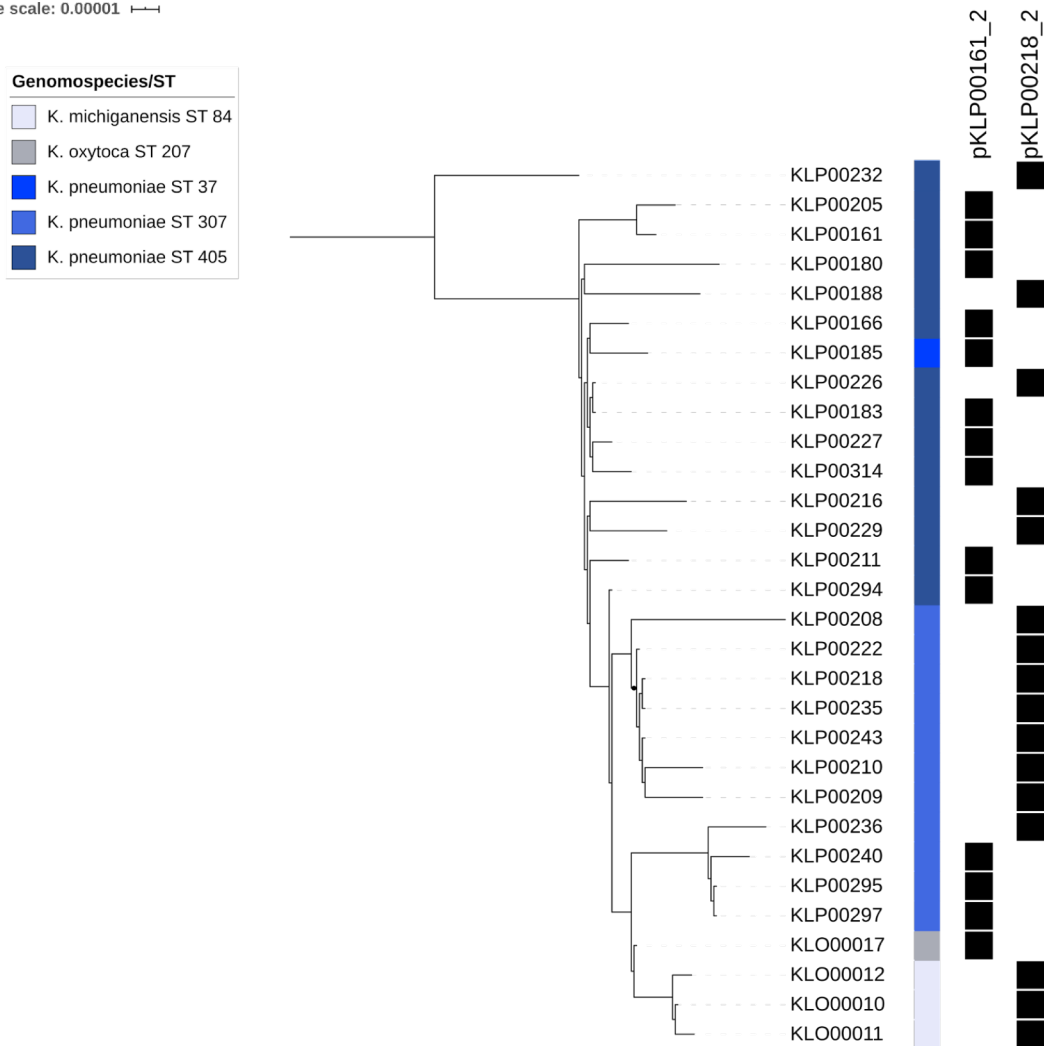
Different bacterial icons denote different species, and different colors denote different sequence types (STs). Figure constructed with icons made by Freepik from Flaticon.com.



Supplementary Figure 2: Alignments of four pairs of plasmids with shared backbone sequences that met sequence identity and coverage thresholds

Grey shading indicates a BLASTn alignment of at least 2,000bp with an e-value of 0 at nucleotide sequence identities of at least 99.9%. Darker shading indicates higher nucleotide sequence identity.

Tree scale: 0.00001



Supplementary Figure 3: Two recombining variants of *Klebsiella* spp. plasmids

Bootstrapped RAxML phylogeny of core plasmid sequences conserved among isolates carrying either variant, using the pKLP00218_2 sequence as a reference.