

1 **Real-time Plasmid Transmission Detection Pipeline**

2 Natalie Scherff^{a,b}#, Jörg Rothgänger^b, Thomas Weniger^b, Alexander Mellmann^a, Dag

3 Harmsen^c

4

5 ^aInstitute of Hygiene, University Hospital Münster, Münster, Germany;

6 ^bRidom GmbH, Münster, Germany;

7 ^cDepartment of Periodontology and Operative Dentistry, University Hospital Münster,

8 Münster, Germany

9

10 Running Head: Plasmid Transmission Detection Pipeline

11

12 #Address correspondence to Natalie Scherff, Natalie.Scherff@ukmuenster.de.

13 <https://orcid.org/0000-0002-8414-0259>

14

15

16 **Abstract**

17 The spread of antimicrobial resistance among bacteria by horizontal plasmid
18 transmissions poses a major challenge for clinical microbiology. Here, we evaluate a new
19 real-time plasmid transmission detection pipeline implemented in the SeqSphere⁺ (Ridom
20 GmbH, Münster, Germany) software.

21 Within the pipeline, a local Mash plasmid database is created and Mash searches with a
22 distance threshold of 0.001 are used to trigger plasmid transmission early warning alerts
23 (EWA). Clonal transmissions are detected using cgMLST allelic differences. The
24 integrated tools MOB-suite, NCBI AMRFinderPlus, CGE MobileElementFinder,
25 pyGenomeViz, and MUMmer are used to characterize plasmids and for visual pairwise
26 plasmid comparisons, respectively. We evaluated the pipeline using published hybrid
27 assemblies (Oxford Nanopore Technology/Illumina) of a surveillance and outbreak
28 dataset with plasmid transmissions. To emulate prospective usage, samples were
29 imported in chronological order of sampling date. Different combinations of the user-
30 adjustable parameters sketch size (1,000 vs 10,000) and plasmid size correction were
31 tested and discrepancies between resulting clusters were analyzed with Quast.

32 When using a sketch size of 1,000 with size correction turned on, the SeqSphere⁺ pipeline
33 agreed with the published data and produced the same clonal and carbapenemase-
34 carrying plasmid clusters. EWAs were in the correct chronological order.

35 In summary, the developed pipeline presented here is suitable for integration into clinical
36 microbiology settings with limited bioinformatics knowledge due to its automated
37 analyses and alert system, which are combined with the GUI-based SeqSphere⁺ platform.
38 Thus, with its integrated sample database, (near) real-time plasmid transmission detection

39 is within reach in bacterial routine-diagnostic settings when long-read sequencing is
40 employed.

41

42 **Importance**

43 Plasmid-mediated spread of antimicrobial resistance (AMR) is a major challenge for
44 clinical microbiology and monitoring of potential plasmid transmissions is essential to
45 combat further dissemination. Whole-genome sequencing (WGS) is often used to surveil
46 nosocomial transmissions but usually limited to the detection of clonal transmissions
47 (based on chromosomal markers). Recent advances in long-read sequencing technologies
48 enable full reconstruction of plasmids and the detection of very similar plasmids but so
49 far easy-to-use bioinformatic tools for this purpose were missing. Here we present an
50 evaluation of an innovative real-time plasmid transmission detection pipeline. It is
51 integrated into the GUI-based SeqSphere⁺ software, which already offers cgMLST based
52 pathogen outbreak detection. It requires very limited bioinformatics knowledge, and its
53 database, automated analyses, and alert system make it well suited for prospective
54 clinical application.

55

56 **Introduction**

57 Antimicrobial resistance (AMR) is a major challenge in clinical microbiology as it is
58 limiting the therapeutic options to manage infectious diseases (1). This is largely
59 mediated by the rapid spread of resistance genes among bacterial populations via
60 horizontal gene transfer (HGT) facilitated by mobile genetic elements (MGE) such as

61 plasmids (2). Plasmids are extrachromosomal DNA elements that can serve as a vehicle
62 to exchange genes between different bacterial strains and even species (3).
63 Recently, multi-species hospital outbreaks of carbapenem-resistant bacteria have been
64 reported (4, 5). Whole-genome sequencing (WGS) based surveillance has become an
65 important tool to monitor nosocomial transmissions in many hospitals (6), but these
66 activities are usually limited to the detection of clonal transmissions, i.e., based on
67 chromosomal markers, of single species. Thus, tools are needed that enable the detection
68 of similar plasmids from WGS data. A variety of software tools has been developed for
69 plasmid reconstruction and characterization (reviewed in (7)). Among the most popular
70 are currently CGE PlasmidFinder (8), COPLA (9), PLACNET (10), and MOB-suite (11),
71 which relies on Mash (12) for comparing k-mers of plasmids. All these tools require
72 some form of bioinformatic expertise, which is often lacking in routine microbiology
73 laboratories. Another major limitation is that plasmids are generally difficult to fully
74 reconstruct from short-reads due to their high number of repetitive regions (7). However,
75 with the recent advances in long-read sequencing, e.g., Oxford Nanopore Technology
76 (ONT), accurate long-read data will become more available to clinical microbiology
77 laboratories (13).
78 The GUI-based Ridom SeqSphere⁺ software (14) is widely used in clinical microbiology
79 laboratories. One of its main purposes is the detection of clonal transmissions using core-
80 genome multi-locus sequence typing (cgMLST) based similarity of bacterial isolates.
81 Here, we evaluated the newly implemented real-time plasmid detection pipeline to enable
82 users to detect potential plasmid transmissions within their own WGS datasets
83 complementary to clonal transmissions. We used two already published datasets of

84 hybrid assemblies of carbapenem-resistant bacteria from a United Kingdom (UK) (15)
85 and a United States (US) hospital (16), respectively.

86

87 **Materials and Methods**

88 SeqSphere⁺ real-time plasmid transmission detection pipeline

89 We developed the pipeline as an extra module for the SeqSphere⁺ (Ridom GmbH,
90 Münster, Germany) software. The first steps involve the reconstruction and
91 characterization of plasmids from assembled data. These steps are based on the MOB-
92 suite (v3.18) software package (11) and work for both short- and long-read data. Short-
93 read and non-circular long read data go into the MOB-recon module first to identify and
94 reconstruct plasmids. The reconstructed plasmids are then forwarded to MOB-typer,
95 which assigns a primary and secondary cluster ID (Mash distance cluster threshold =
96 0.006 and 0.025, respectively), rep and relaxase type(s), and mobility prediction to each
97 plasmid. For long-read data, the FASTA headers are checked for circularity terms
98 ([topology = circular]) and only non-circular contigs are forwarded to MOB-recon. Circular
99 contigs that are above 500kb in size are considered to be chromosomes, smaller contigs
100 are considered to be plasmids and directly go into MOB-typer. Optionally, the MOB-
101 recon module can also be skipped using a [non-recon] term in the FASTA header. AMR
102 targets are determined using the NCBI AMRFinderPlus (v3.11.26) (17). Integrated
103 mobile genetic elements (iMGEs) are detected with CGE MobileElementFinder (v1.1.2)
104 (18). The typing results from MOB-typer, AMRFinderPlus, and MobileElementFinder
105 are then summarized in a tabular “Chromosome & Plasmid” overview. From here,
106 plasmid contigs can be exported as FASTA files for downstream analyses.

107 AMRFinderPlus results are grouped into two columns, one of them only showing priority
108 AMR genes, which are defined as targets that might confer resistance to carbapenems,
109 colistin, vancomycin, or methicillin, and genes that contain ESBL or AmpC in their
110 name. The detected iMGE(s) are also shown in a separate column if they enclose a
111 priority AMR gene.

112 Long-read data plasmid transmission analysis module

113 With the recently introduced version 10 of the SeqSphere⁺ software with the long-read
114 data plasmid transmission analysis module, potential plasmid transmissions can be
115 detected using Mash (v2.1) comparisons (12). The identified plasmids can be used to
116 build a local Mash plasmid database from the users' own data. Pairwise Mash distances
117 can then be utilized to detect very similar plasmids. This information can be used in
118 several ways. First, if new plasmids are imported via the pipeline mode, they are
119 compared with all existing plasmids in the database and if there are matches below a
120 certain threshold (default = 0.001), an early warning alert (EWA) is triggered and shown
121 on the SeqSphere⁺ start screen. No EWA is generated if the respective samples contain a
122 cgMLST scheme and have a clonal (i.e., cgMLST allelic distance below the species-
123 specific threshold) relationship. Second, the complete plasmid database can be used to
124 perform an all-against-all comparison to create an exportable distance matrix and/or a
125 single-linkage clustering. Third, the database can be searched for similar plasmids from a
126 single plasmid of interest. This plasmid search can be triggered from both the
127 "Chromosome & Plasmid" overview and EWA reports. Within the resulting table,
128 sample metadata and plasmid typing results for each plasmid are included. Finally,
129 partially annotated plasmids (origin(s) of replication, priority AMR gene(s), other AMR

130 gene(s), and iMGE(s)) can also be pairwise aligned and visualized using MUMmer v3.23
131 (19) and pyGenomeViz (v0.4.4; <https://moshi4.github.io/pyGenomeViz/>). To create a
132 better visualization, there is an optional fixstart function to change the start and
133 orientation of contigs. It uses *dnaA* for chromosomes and for plasmids the first
134 occurrence of a *rep* gene as the start of the reoriented contig. Figure 1 shows an overview
135 of the complete pipeline.

136 SeqSphere⁺ allows for configuration of multiple parameters including sketch size,
137 sequence length compensation, and Mash distance threshold for matches (default:0.001).
138 For the sketch size the user has the choice between 1,000 and 10,000. The k-mer size of
139 21 is fixed. Due to the large variability in plasmid length and as Mash distances favor
140 contigs of similar size we have implemented a scaling factor for comparing plasmids with
141 difference in unique content where the match threshold is lowered for each 1% of size
142 difference. The default is a lowering of 0.0003 per 1% size difference up to a limit of
143 40% size difference but both the lowering value and the upper limit can be adjusted by
144 the user.

145

146 Evaluation of the pipeline using two published datasets

147 To test the pipeline functionality and evaluate parameter settings, the pipeline was
148 applied to two different published datasets. The first was a set of 85 isolates of
149 phenotypically carbapenem-resistant Gram-negative bacteria from a single hospital in the
150 UK (Addenbrooke's Hospital, Cambridge) that were collected over a period of six years
151 (Addenbrookes dataset) (15). We excluded four samples without a detected
152 carbapenemase gene leading to a total number of 81 samples in our test dataset. These

153 isolates were originally sequenced using a combination of Illumina HiSeq and Oxford
154 Nanopore Technologies (ONT) MinION R9.4.1 technologies. Original hybrid assemblies
155 were created using Flye (v2.8) (20), Unicycler (v0.4.8) (21) and Canu (v2.1.1) (22).
156 These hybrid assemblies were downloaded in FASTA format from the European
157 Nucleotide Archive (ENA), where they had been previously deposited under the study
158 accession PRJEB30134. A Mash search against the NCBI RefSeq genomes was
159 performed to confirm species assignment and rule out potential contamination. From the
160 information in the original publication, we did assume that each contig, except the
161 chromosome, was a single plasmid. However, since we did not have circularity
162 information, we used a [non-recon] term in the FASTA headers to skip the MOB-recon
163 part as additional plasmid reconstruction was not done in the original publication. The
164 study was done retrospectively but to test the EWAs, we emulated a prospective analysis
165 by importing the files in the chronological order of sampling date. Sampling date and
166 patient information were taken from the supplemental material provided with the original
167 publication.

168 Plasmid clusters were defined as plasmids with a Mash distance ≤ 0.001 . In addition to
169 clusters, the number of potential transmissions was counted. If a sample carried multiple
170 plasmids that belonged to different plasmid clusters, it was assessed whether all of these
171 plasmids were transferred to the same or different receiving samples. If all plasmids were
172 shared with the same sample, this was counted as one co-transmission event. If plasmids
173 of one sample were connected to different samples, each connection was counted as a
174 single plasmid transmission event. If two samples had a clonal relation, the respective
175 plasmids were not counted as a plasmid but as a clonal transmission event.

176 To assess clonal relations, we used cgMLST allelic distances. For species with public
177 cgMLST schemes hosted on the Ridom nomenclature server (www.cgMLST.org), we
178 used the species-specific default threshold. For species without a public scheme, we
179 created *ad hoc* schemes, using the NCBI Genomes reference sequence of each species as
180 seed genome and a threshold of 15 alleles to determine clonal clusters. To check for
181 additional transmissions of only iMGEs, a user-defined task template was created from
182 an allele database that contained all iMGE with a carbapenemase gene that were found
183 using the CGE MobileElementFinder with thresholds of $\geq 90\%$ identity and $\geq 95\%$
184 alignment length.

185 For parameter evaluation, we tested four combinations: a sketch size of 1,000 and 10,000,
186 each with and without size compensation applied to the Mash distance. A cluster
187 threshold of ≤ 0.001 was used for all four approaches. Plasmid clusters with
188 discrepancies between these different runs were analyzed with Quast (v5.2) (23). Here,
189 the longest plasmid of each cluster was set as reference. Mismatches, InDels, and larger
190 gaps were analyzed and counted. . Finally, the pipeline was applied to the non
191 carbapenemase-carrying plasmids of this dataset as well.

192 To prove that the pipeline with the determined Mash parameters also works in another
193 setting with different strains, we evaluated a second dataset. This one contained 19
194 isolates of a multispecies ($n = 7$) outbreak of *bla*NDM-5-producing Enterobacterales in a
195 US hospital (UPMC Presbyterian dataset) (16). The isolates were collected from
196 February 2021 to February 2023 from 15 patients. All isolates were sequenced on an
197 Illumina NextSeq 550 and ONT MinION with R9.4.1 flow cells. Hybrid assemblies were
198 created using Unicycer v0.5.0 and deposited by the original authors for download at

199 NCBI under BioProject PRJNA981541. In the original paper, outbreak plasmids were
200 defined as *IncX3* plasmids harboring *bla*NDM-5, “that had sequence coverage values
201 (*i.e.*, the proportion of each plasmid’s sequence that was found in the Illumina genome)
202 of at least 95%, and a sequence identity of <15 single nucleotide polymorphisms (SNPs)
203 per 100 kb of plasmid-sequence compared with the fully resolved, circular plasmid
204 sequence” (16). For our evaluation, the same analysis as described above was applied to
205 the carbapenemase-carrying plasmids of this dataset.

206

207 **Results**

208 Analysis of clonal clusters in the Addenbrookes dataset

209 Clonal clusters were determined based on cgMLST allelic distances with species-specific
210 thresholds (15 alleles difference for species without public cgMLST schemes). In total,
211 seven clonal clusters were detected, three for *Klebsiella pneumoniae*, two for *Escherichia*
212 *coli*, one for *Pseudomonas putida*, and one for *Enterobacter hormaechei* (Table 1).
213 Except one *K. pneumoniae* cluster comprising eleven isolates (kleb_2), all clusters
214 contained two isolates. The allelic distances were between 0 and 11. In addition, we
215 found two *E. roggkampii* isolates with a distance of one allele but since these isolates
216 were from the same patient, they were denoted as intra-host variability.

217 Analysis of carbapenemase-carrying plasmids in the Addenbrookes dataset

218 Initially, we only analyzed the plasmids carrying a carbapenemase gene. Out of the 81
219 isolates, eleven isolates contained chromosomally encoded carbapenemase genes. The
220 other 70 isolates had carbapenemase-carrying plasmids. One isolate comprised two
221 carbapenemase-carrying plasmids leading to a total number of 71 analyzed plasmids.

222 By re-sequencing strains from documented and published plasmid transmission events
223 with ONT and Pacific Biosciences sequencing methodology (4, 24), we determined a
224 Mash distance of 0.001 as suitable for recovering nearly identical plasmids (data not
225 shown). By varying the unique content of those plasmids, we came up with a scaling
226 factor of 0.0003 per 1% plasmid size difference. Next, using a sketch size of 10,000, a
227 Mash distance threshold of 0.001, and a 0.0003 size compensation, we found nine
228 plasmid clusters (Table 1). One additional plasmid cluster consisted of only clonal
229 isolates (kleb_2). Three plasmid clusters (Clust_148, Clust_9, Clust_72) partly contained
230 clonal isolates. The clusters comprised between two and eight isolates from one to three
231 different species. The relatedness of the 71 carbapenemase-carrying plasmids is
232 visualized with a population snapshot using the exported clustering distance matrix
233 (Figure 2).

234 To evaluate the Mash parameters, we repeated this analysis with a lower sketch size of
235 1,000 and tested both sketch sizes with and without size compensation and compared the
236 resulting clusters. Out of the nine plasmid clusters, four showed discrepancies between
237 the different methods (Table 2). One plasmid pair (Clust_54) was missed with a sketch
238 size of 1,000 but only when size compensation was turned off. In Clust_9, one additional
239 sample was included in the cluster for both sketch sizes when size compensation was
240 applied. Clust_148 also contained an additional plasmid with size compensation, but only
241 with a sketch size of 10,000. Finally, Clust_98/Clust_217 was lumped together and
242 contained nine samples with size compensation but without size compensation it was
243 divided in two clusters with six and three samples, respectively. We counted the number
244 of differences between the disputable plasmid and the reference plasmid of each cluster.

245 For Clust_54 there were three, for Clust_148 five, and for Clust_9 13 differences,
246 respectively. For Clust_98/Clust_217 we compared all samples of Clust_217 with one
247 reference of Clust_98, which resulted in two differences.

248 In one cluster (Clust_72), we detected a case, where one plasmid (size = 12,282 bp)
249 contained a duplicated sequence, which was otherwise completely identical with another
250 plasmid (size = 6,141 bp) that only contained the single sequence, i.e. their uncorrected
251 Mash distance was 0. We checked for additional transmissions of iMGEs by retrieving
252 the sequences of all iMGEs containing a carbapenemase gene and searching for these
253 sequences in all isolates utilizing a user-defined task template. No additional
254 transmissions of iMGEs were found using this approach.

255 Analysis of all plasmids in the Addenbrookes dataset

256 After analysis of the carbapenemase-carrying plasmids, in contrast to the original
257 publication, we also had a look at the overall plasmid landscape of these samples. In total,
258 324 plasmids were found in 78 of the 81 samples (three samples did not contain
259 plasmids). Each sample carried 1-10 plasmids with an average of four plasmids per
260 sample. Of these, 175 plasmids from 58 samples belonged to a cluster (Mash distance
261 0.001). On average, three plasmids per isolate belonged to a cluster with more than one
262 plasmid. A total of 51 clusters (including clonal and carbapenemase-carrying plasmid
263 clusters) was detected with 2-11 isolates per cluster (median = 2, mean = 3.4).

264 In addition to clonal and plasmid clusters, we counted the number of potential
265 transmission events. In total, 44 single transmission events were counted. Of these,
266 twelve were clonal transmissions, 18 were single plasmid transmissions, and 14 were co-
267 transmissions of two or more plasmids to the same sample. Three of the single and five of

268 the co-plasmid-transmissions were intra-host transmissions. Fifteen isolates transferred
269 more than one plasmid; for nine of them, all plasmids were transferred to the same and
270 for six of them to different receiving samples. Furthermore, we detected two cases where
271 an isolate received plasmids from a clonal and an additional plasmid transmission. One of
272 these involved two different species, *K. pneumoniae* and *E. coli* (Figure 3). Here, the two
273 *K. pneumoniae* samples had identical MLST sequence types (STs) and cgMLST allelic
274 profiles. However, the sample that was isolated later carried additional plasmids. Two of
275 these plasmids (size = 1.8 kb & 5.2 kb) clustered with plasmids of an *E. coli* isolate that
276 also shared a 145 kb *bla*NDM1-carrying plasmid with both *K. pneumoniae* samples.

277 Analysis of carbapenemase-carrying plasmids in the UPMC Presbyterian dataset

278 As a further evaluation, we analyzed the carbapenemase-carrying plasmids of a second
279 dataset, comprising 19 isolates of a multispecies *bla*NDM-5 outbreak (16). According to
280 the original publication, all isolates belonged to a single outbreak, where *bla*NDM-5 was
281 encoded on a 46.2 kb *IncX3* family plasmid. The plasmids differed by only 0-2 SNPs
282 with the exception of one shorter (45.2 kb) and one larger (55.8 kb) plasmid that lost or
283 acquired a transposase gene, respectively.

284 Using our pipeline with a sketch size of 1,000 and size correction turned on, all 19
285 plasmids were put into a single plasmid cluster. However, when size correction was
286 turned off, the larger (17% size difference) plasmid did not cluster with the remaining
287 plasmids (uncorrected Mash distance = 0.003). Both, *bla*NDM-5 and *IncX3* were detected
288 on all plasmids by AMRFinderPlus and MOB-suite, respectively. In addition, among
289 other information, relaxase type MOB_P was assigned by MOB-typer and the plasmids
290 were predicted to be conjugative.

291

292 **Discussion**

293 In this study, we used two published datasets of plasmids from carbapenem-resistant
294 bacteria to evaluate a new plasmid detection pipeline implemented in the SeqSphere⁺
295 software. We determined a Mash distance of 0.001 and a scaling factor of 0.0003 per 1%
296 plasmid size difference as suitable for recovering nearly identical plasmids. Thresholds
297 like these are always somewhat arbitrary and need to be a compromise between
298 sensitivity and specificity. However, the developers of the MOB-suite came up with a
299 distance threshold of 0.025 for short-read sequencing data (25). With long-read data this
300 value can be lower. Independent from our own considerations, Roberts *et al.* came up
301 with the same threshold of 0.001 (15). For both datasets, our results are in complete
302 agreement with the findings from the original studies, in particular, the same clonal and
303 plasmid clusters were found. Moreover, the same resistance genes and plasmid rep-types
304 were detected, despite using different tools and databases. In the Roberts *et al.* study,
305 Abricate (<https://github.com/tseemann/abricate>) was used with the CARD (26) and
306 PlasmidFinder (8) databases, while our pipeline is based on NCBI AMRFinderPlus (17)
307 and MOB-suite (11). In both the Roberts *et al.* and our study, Mash with a sketch size of
308 10,000 and a distance threshold of 0.001 was used to detect plasmid results, leading not
309 surprisingly to the same cluster results. However, in contrast to a file-based Mash
310 approach, our automated, database-based pipeline enables a prospective real-time
311 surveillance of plasmid transmissions.

312 Raabe *et al.* defined outbreak plasmids based on Illumina sequence coverage values of at
313 least 95%, and a sequence identity of <15 SNPs per 100 kb of plasmid sequence when

314 compared against the complete reference plasmid sequence. They also used the CGE
315 PlasmidFinder for rep-typing. Different approaches were also used for clonal cluster
316 detection. While Roberts *et al.* used a combination of SNPs and split k-mer analysis
317 (SKA), the SeqSphere⁺ pipeline is based on cgMLST allelic distances.

318 .

319 We compared the clustering results of Mash databases constructed with a sketch size of
320 1,000 and 10,000. Both are values used for plasmids in the literature (15, 25). While the
321 higher sketch size increases the sensitivity, it also requires more computational memory
322 and power and enlarges the database size. We found that a sketch size of 1,000 was
323 generally enough to detect the same clusters at a threshold of 0.001, with one exception.
324 The sensitivity can be increased without a significant increase in computational power
325 requirements by applying a size compensation. The size compensation allows to detect
326 potential transmissions of plasmids that share the same backbone but may have lost or
327 gained larger unique segments. This was shown with the UPMC Presbyterian dataset,
328 where one insertion of a transposase gene increased the size of the otherwise identical
329 plasmid by 17 %. The Mash distance between this larger plasmid and the other outbreak
330 plasmids was above the clustering threshold without size correction. However, size
331 compensation must have an upper limit to prevent too many false positives just resulting
332 from shared promiscuous iMGes. An interesting finding was a plasmid pair, where one
333 plasmid consisted of two exact copies of another plasmid. This is likely a result of
334 multimer formation. To account for such phenomena, we decided to limit the size
335 compensation to a maximum of 40% plasmid size differences but keep the uncorrected
336 Mash distance value for plasmids with a larger size difference instead of ignoring

337 plasmid pairs above 40% size difference. Thus, plasmids with large size differences but
338 very similar content, as for example multimers, can still be detected.

339 When analyzing the whole plasmid landscape of the Addenbrookes dataset, we found that
340 the reconstruction of potential transmission scenarios can be complicated due to complex
341 co-transmission patterns. Based on genetic clustering, the number of plasmid
342 transmission events within the analyzed dataset is estimated to be twice as common as
343 clonal transmission of pathogens but detailed epidemiological information is needed to
344 confirm these results. In particular, it is still largely unknown how often very similar
345 plasmids occur without a transmission, i.e. as “core plasmids” of certain lineages or
346 species. As other plasmid comparison tools, the pipeline only works with well
347 reconstructed, ideally circular plasmids and thus requires long-read sequences. Moreover,
348 Mash as a k-mer-based approach does not take synteny into account. Although same
349 content is usually also same synteny in a local context, it is possible to check for
350 rearrangements between plasmid sequences by utilizing pyGenomeViz, which is
351 implemented in SeqSphere+, as a visual approach. For more diverse data, alternatively, a
352 new tool such as Pling (27), which is alignment-based and thereby more compute-
353 intensive, could be used. Currently the literature is still lacking analyses of plasmid
354 transmissions based on long-read data. However, with long-read sequencing becoming
355 more and more achievable for clinical laboratories, the need for suitable analysis tools is
356 growing. We presented here a retrospective re-analysis of published datasets. However,
357 in a prospective surveillance scenario, every single transmission would be captured by
358 the EWA system as soon as the isolates are imported, allowing for a near real-time
359 plasmid transmission detection.

360 In conclusion, the new SeqSphere⁺ plasmid transmission detection pipeline allows the
361 detection of plasmid clusters in both retrospective (single linkage clustering of plasmid
362 databases) and prospective (Early Warning Alerts) studies. The automated, database-
363 based approach enables near real-time surveillance and allows for rapid informed actions.
364

365 **References**

- 366 1. Murray CJ, Ikuta KS, Sharara F, Swetschinski L, Aguilar GR, Gray A, Han C,
367 Bisignano C, Rao P, Wool E, Johnson SC, Browne AJ, Chipeta MG, Fell F, Hackett
368 S, Haines-Woodhouse G, Hamadani BHK, Kumaran EAP, McManigal B, Agarwal
369 R, Akech S, Albertson S, Amuasi J, Andrews J, Aravkin A, Ashley E, Bailey F,
370 Baker S, Basnyat B, Bekker A, Bender R, Bethou A, Bielicki J, Boonkasidecha S,
371 Bukosia J, Carneiro C, Castañeda-Orjuela C, Chansamouth V, Chaurasia S,
372 Chiurchiù S, Chowdhury F, Cook AJ, Cooper B, Cressey TR, Criollo-Mora E,
373 Cunningham M, Darboe S, Day NPJ, Luca MD, Dokova K, Dramowski A, Dunachie
374 SJ, Eckmanns T, Eibach D, Emami A, Feasey N, Fisher-Pearson N, Forrest K,
375 Garrett D, Gastmeier P, Giref AZ, Greer RC, Gupta V, Haller S, Haselbeck A, Hay
376 SI, Holm M, Hopkins S, Iregbu KC, Jacobs J, Jarovsky D, Javanmardi F, Khorana M,
377 Kisson N, Kobeissi E, Kostyanov T, Krapp F, Krumkamp R, Kumar A, Kyu HH,
378 Lim C, Limmathurotsakul D, Loftus MJ, Lunn M, Ma J, Mturi N, Munera-Huertas T,
379 Musicha P, Mussi-Pinhata MM, Nakamura T, Nanavati R, Nangia S, Newton P,
380 Ngoun C, Novotney A, Nwakanma D, Obiero CW, Olivares-Martinez A, Olliaro P,
381 Ooko E, Ortiz-Brizuela E, Peleg AY, Perrone C, Plakkal N, Ponce-de-Leon A, Raad
382 M, Ramdin T, Riddell A, Roberts T, Robotham JV, Roca A, Rudd KE, Russell N,
383 Schnall J, Scott JAG, Shivamallappa M, Sifuentes-Osornio J, Steenkeste N,
384 Stewardson AJ, Stoeva T, Tasak N, Thaiprakong A, Thwaites G, Turner C, Turner P,
385 Doorn HR van, Velaphi S, Vongpradith A, Vu H, Walsh T, Waner S,
386 Wangrangsimakul T, Wozniak T, Zheng P, Sartorius B, Lopez AD, Stergachis A,

- 387 Moore C, Dolecek C, Naghavi M. 2022. Global burden of bacterial antimicrobial
388 resistance in 2019: a systematic analysis. *The Lancet* 399:629–655.
- 389 2. Sheppard AE, Stoesser N, Wilson DJ, Sebra R, Kasarskis A, Anson LW, Giess A,
390 Pankhurst LJ, Vaughan A, Grim CJ, Cox HL, Yeh AJ, Modernising Medical
391 Microbiology (MMM) Informatics Group, Sifri CD, Walker AS, Peto TE, Crook
392 DW, Mathers AJ. 2016. Nested Russian doll-like genetic mobility drives rapid
393 dissemination of the carbapenem resistance gene *blaKPC*. *Antimicrob Agents*
394 *Chemother* 60:3767–3778.
- 395 3. Orlek A, Stoesser N, Anjum MF, Doumith M, Ellington MJ, Peto T, Crook D,
396 Woodford N, Sarah Walker A, Phan H, Sheppard AE. 2017. Plasmid classification in
397 an era of whole-genome sequencing: Application in studies of antibiotic resistance
398 epidemiology. *Frontiers in Microbiology* 8:1–10.
- 399 4. Weber RE, Pietsch M, Frühauf A, Pfeifer Y, Martin M, Luft D, Gatermann S,
400 Pfennigwerth N, Kaase M, Werner G, Fuchs S. 2019. IS26-mediated transfer of
401 *blaNDM-1* as the main route of resistance transmission during a polyclonal,
402 multispecies outbreak in a German hospital. *Frontiers in Microbiology* 10:1–14.
- 403 5. Marimuthu K, Venkatachalam I, Koh V, Harbarth S, Perencevich E, Cherng BPZ,
404 Fong RKC, Pada SK, Ooi ST, Smitasin N, Thoon KC, Tambyah PA, Hsu LY, Koh
405 TH, De PP, Tan TY, Chan D, Deepak RN, Tee NWS, Kwa A, Cai Y, Teo Y-Y,
406 Thevasagayam NM, Prakki SRS, Xu W, Khong WX, Henderson D, Stoesser N, Eyre
407 DW, Crook D, Ang M, Lin RTP, Chow A, Cook AR, Teo J, Ng OT, Carbapenemase-
408 Producing Enterobacteriaceae in Singapore (CaPES) Study Group, Marimuthu K,

- 409 Venkatachalam I, Cherng BPZ, Fong RKC, Pada SK, Ooi ST, Smitasin N, Thoon
410 KC, Hsu LY, Koh TH, De PP, Tan TY, Chan D, Deepak RN, Tee NWS, Ang M, Lin
411 RTP, Teo J, Ng OT. 2022. Whole genome sequencing reveals hidden transmission of
412 carbapenemase-producing Enterobacterales. *Nat Commun* 13:3052.
- 413 6. Mellmann A, Bletz S, Böking T, Kipp F, Becker K, Schultes A, Prior K, Harmsen D.
414 2016. Real-time genome sequencing of resistant bacteria provides precision infection
415 control in an institutional setting. *Journal of Clinical Microbiology* 54:2874–2881.
- 416 7. Paganini JA, Plantinga NL, Arredondo-Alonso S, Willems RJL, Schürch AC. 2021.
417 Recovering *Escherichia coli* plasmids in the absence of long-read sequencing data.
418 *Microorganisms* 9:1613.
- 419 8. Carattoli A, Zankari E, Garcia-Fernandez A, Voldby Larsen M, Lund O, Villa L,
420 Moller Aarestrup F, Hasman H. 2014. *In silico* detection and typing of plasmids
421 using PlasmidFinder and plasmid multilocus sequence typing. *Antimicrobial agents
422 and chemotherapy* 58:3895–3903.
- 423 9. Redondo-Salvo S, Bartomeus-Peñalver R, Vielva L, Tagg KA, Webb HE, Fernández-
424 López R, de la Cruz F. 2021. COPLA, a taxonomic classifier of plasmids. *BMC
425 Bioinformatics* 22:390.
- 426 10. Vielva L, de Toro M, Lanza VF, de la Cruz F. 2017. PLACNETw: a web-based tool
427 for plasmid reconstruction from bacterial genomes. *Bioinformatics* 33:3796–3798.
- 428 11. Robertson J, Nash JHE. 2018. MOB-suite: software tools for clustering,
429 reconstruction and typing of plasmids from draft assemblies. *Microbial Genomics* 4.

- 430 12. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy
431 AM. 2016. Mash: Fast genome and metagenome distance estimation using MinHash.
432 *Genome Biology* 17:1–14.
- 433 13. Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA, Wollenberg
434 RD, Albertsen M. 2022. Oxford Nanopore R10.4 long-read sequencing enables the
435 generation of near-finished bacterial genomes from pure cultures and metagenomes
436 without short-read or reference polishing. *Nat Methods* 19:823–826.
- 437 14. Jünemann S, Sedlazeck FJ, Prior K, Albersmeier A, John U, Kalinowski J, Mellmann
438 A, Goesmann A, von Haeseler A, Stoye J, Harmsen D. 2013. Updating benchtop
439 sequencing performance comparison. *Nature biotechnology* 31:294–296.
- 440 15. Roberts LW, Enoch DA, Khokhar F, Blackwell GA, Wilson H, Warne B, Gouliouris
441 T, Iqbal Z, Török ME. 2023. Long-read sequencing reveals genomic diversity and
442 associated plasmid movement of carbapenemase-producing bacteria in a UK hospital
443 over 6 years. *Microbial Genomics* 9:001048.
- 444 16. Raabe NJ, Valek AL, Griffith MP, Mills E, Waggle K, Srinivasa VR, Ayres AM,
445 Bradford C, Creager HM, Pless LL, Sundermann AJ, Van Tyne D, Snyder GM,
446 Harrison LH. 2024. Real-time genomic epidemiologic investigation of a multispecies
447 plasmid-associated hospital outbreak of NDM-5-producing Enterobacterales
448 infections. *International Journal of Infectious Diseases* 142:106971.
- 449 17. Feldgarden M, Brover V, Gonzalez-Escalona N, Frye JG, Haendiges J, Haft DH,
450 Hoffmann M, Pettengill JB, Prasad AB, Tillman GE, Tyson GH, Klimke W. 2021.

- 451 AMRFinderPlus and the Reference Gene Catalog facilitate examination of the
452 genomic links among antimicrobial resistance, stress response, and virulence.
453 Scientific Reports 11:1–9.
- 454 18. Johansson MHK, Bortolaia V, Tansirichaiya S, Aarestrup FM, Roberts AP, Petersen
455 TN. 2021. Detection of mobile genetic elements associated with antibiotic resistance
456 in *Salmonella enterica* using a newly developed web tool: MobileElementFinder. J
457 Antimicrob Chemother 76:101–109.
- 458 19. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg
459 SL. 2004. Versatile and open software for comparing large genomes. Genome Biol
460 5:R12.
- 461 20. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone
462 reads using repeat graphs. Nat Biotechnol 37:540–546.
- 463 21. Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Unicycler: Resolving bacterial
464 genome assemblies from short and long sequencing reads. PLoS Comput Biol
465 13:e1005595.
- 466 22. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu:
467 scalable and accurate long-read assembly via adaptive k-mer weighting and repeat
468 separation. Genome Res 27:722–736.
- 469 23. Mikheenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A. 2018. Versatile
470 genome assembly evaluation with QUAST-LG. Bioinformatics 34:i142–i150.

- 471 24. van Almsick V, Schuler F, Mellmann A, Schwierzeck V. 2022. The use of long-read
472 sequencing technologies in infection control: Horizontal transfer of a *bla*CTX-M-27
473 containing IncFII plasmid in a patient screening sample. *Microorganisms* 10:491.
- 474 25. Robertson J, Bessonov K, Schonfeld J, Nash JHE. 2020. Universal whole-sequence-
475 based plasmid typing and its utility to prediction of host range and epidemiological
476 surveillance. *Microbial Genomics* 6:1–12.
- 477 26. Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P, Tsang KK, Lago BA, Dave
478 BM, Pereira S, Sharma AN, Doshi S, Courtot MM, Lo R, Williams LE, Frye JG,
479 Elsayegh T, Sardar D, Westman EL, Pawlowski AC, Johnson TA, Brinkman FSL,
480 Wright GD, McArthur AG. 2017. CARD 2017: Expansion and model-centric
481 curation of the comprehensive antibiotic resistance database. *Nucleic Acids Research*
482 45:D566–D573.
- 483 27. Frolova D, Lima L, Roberts L, Bohnenkämper L, Wittler R, Stoye J, Iqbal Z. 2024.
484 Applying rearrangement distances to enable plasmid epidemiology with pling.
485 bioRxiv <https://doi.org/10.1101/2024.06.12.598623>.

486

487

488 **Acknowledgments**

489 We thank Leah Roberts (CIIC, Queensland University of Technology) for providing help
490 with sequence and epidemiological data. Further, we thank James Robertson and Kyrilo
491 Bessonov (both National Microbiology Laboratory, Public Health Agency of Canada) for
492 critical reading and improvement of the manuscript.

493

494 **Disclaimers**

495 NS, JR, and TW are (part-time) employees of Ridom GmbH. JR, TW, and DH are
496 shareholders of Ridom GmbH. AM declares no conflict of interest.

497

498

499

500

501 **Tables**

502 **Table 1:** Overview of detected clonal and carbapenemase-carrying plasmid transmission
 503 clusters. Cluster names correspond to the original publication.

504

| Cluster | No. of isolates | Species | Carbapenemase gene | Rep types |
|-------------------------------|-----------------|--|--------------------|---------------|
| Clonal clusters | | | | |
| pseu_1 | 2 | <i>Pseudomonas putida</i> | IMP-70 | |
| ecoli_1 | 2 | <i>Escherichia coli</i> | OXA-232 | |
| ecoli_2 | 2 | <i>Escherichia coli</i> | NDM-5 | |
| ehor_1 | 2 | <i>Enterobacter hormaechei</i> | IMP-70 | |
| kleb_1 | 2 | <i>Klebsiella pneumoniae</i> | NDM-1 | |
| kleb_2 | 11 | <i>Klebsiella pneumoniae</i> | NDM-1 | |
| kleb_3 | 2 | <i>Klebsiella pneumoniae</i> | OXA-232 | |
| Intra-host variability | | | | |
| erog_1 | 2 | <i>Enterobacter roggenkampii</i> | OXA-181 | |
| Plasmid clusters | | | | |
| Clust_33 | 2 | <i>Escherichia coli, Klebsiella pneumoniae</i> | NDM-6 | IncFIB |
| Clust_133 | 2 | <i>Acinetobacter baumannii</i> | OXA-72 | rep_1127 |
| Clust_17 | 4 | <i>Escherichia coli, Klebsiella aerogenes, Klebsiella pneumoniae</i> | OXA-48 | IncI |
| Clust_148* | 4 | <i>Escherichia coli, Klebsiella pneumoniae</i> | NDM-1 | IncC |
| Clust_54 | 2 | <i>Escherichia coli</i> | NDM-5 | IncFIB/IncFIA |
| Clust_9* | 3 | <i>Enterobacter hormaechei</i> | IMP-70 | IncC |
| Clust_72* | 8 | <i>Escherichia coli, Klebsiella pneumoniae</i> | OXA-232 | Col |
| Clust_98 | 6 | <i>Escherichia coli, Enterobacter hormaechei, Klebsiella pneumoniae</i> | OXA-181 | IncX3 |
| Clust_217 | 4 | <i>Enterobacter roggenkampii, Klebsiella grimontii, Raoultella ornithinolytica</i> | OXA-181 | IncX3 |

* includes clonal isolates

505 **Table 2:** Plasmids with variations in clustering results between different Mash database
 506 parameters. (Compensated) Mash distances and Quast results for each comparison of a
 507 plasmid with a reference plasmid (grey fill color) are shown. The different parameters
 508 were sketch size of 1,000 (1k) or 10,000 (10k) and use of size compensation (sc). Mash
 509 distance values that led to different clustering results (threshold = 0.001) are highlighted
 510 in orange. The dashed line marks two clusters that were joined into one cluster with size
 511 correction but separated without.
 512

| Plasmid clusters | | | Mash distance to reference | | | | Quast results | | | |
|------------------|-------------------|--------------------------------|----------------------------|---------|---------|----------|---------------|--------|------------|------------|
| Cluster | Plasmid size (kb) | Size difference to reference % | 1k | 1k + sc | 10k | 10k + sc | Mismatches | Indels | Large gaps | Sum events |
| Clust_9 | 127,2 | 0% | | | | | | | | |
| | 127,2 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | 0 | 0 | 0 | 0 |
| | 126,2 | 1% | 0,0002 | 0,0000 | 0,0000 | 0,0000 | 0 | 1 | 2 | 3 |
| | 96,3 | 24% | 0,0068 | 0,0000 | 0,0074 | 0,0002 | 1 | 10 | 2 | 13 |
| Clust_148 | 145,8 | 0% | | | | | | | | |
| | 145,8 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | 0 | 0 | 0 | 0 |
| | 145,8 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | 0 | 0 | 0 | 0 |
| | 140,1 | 4% | 0,0022 | 0,00103 | 0,0019 | 0,0007 | 4 | 0 | 1 | 5 |
| Clust_54 | 153,2 | 0% | | | | | | | | |
| | 148,4 | 3% | 0,00101 | 0,0001 | 0,00099 | 0,00004 | 1 | 1 | 1 | 3 |
| Clust_98 | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| | 51,5 | 0% | 0,0000 | 0,0000 | 0,0000 | 0,0000 | | | | |
| Clust_217 | 47,1 | 9% | 0,0021 | 0,0000 | 0,0022 | 0 | 1 | 0 | 1 | 2 |
| | 47,1 | 9% | 0,0021 | 0,0000 | 0,0022 | 0 | 1 | 0 | 1 | 2 |
| | 48,5 | 6% | 0,0022 | 0,0004 | 0,0022 | 0,0005 | 1 | 0 | 1 | 2 |

513

514

515 **Figures**

516

517 **Figure 1:** Overview of plasmid transmission detection pipeline.

518

519 **Figure 2:** Plasmid population snapshot based on Mash distances between 71
520 carbapenemase-carrying plasmids. Each circle represents one plasmid, plasmids with a
521 Mash distance of 0 are lumped together in single circles with a larger size. Connection
522 lines are drawn for plasmids with a Mash distance <0.06 . Fill color represents the species
523 of the respective isolates and colored, dotted circles mark clonal clusters among these
524 isolates. Plasmid clusters (Mash distance threshold <0.001) are illustrated by grey shaded
525 areas around the circles. Clonal (in color) and plasmid (in grey) cluster names are stated
526 next to the clusters.

527

528

529 **Figure 3:** Visualization of a combination of clonal and horizontal plasmid transmissions.
530 Direction of transmission is assumed based on temporal order of isolation. ST = MLST
531 sequence type.

Assembled Illumina Data

Non-circular LR Data

Assembled Circular LR* Data

(optional Fix rep/dnaA Start and Orientation)

MOB-recon

(includes MOB-typer)

MOB-typer

Plasmid Overview

(with MOB-suite, NCBI AMR Finder, and CGE MobileElement Finder)

Plasmid Reconstruction & Characterization

own plasmids

Long-read Data Plasmid Transmission Analysis Module

**Mash
Plasmid
DB(s)**

optional SLC**

**Plasmid
Transmission EWA*****

plasmid size difference compensation,
spatiotemporal limiters (excluding likely
clonal events)

Plasmid Transmission

**Plasmid
Search Report**

plasmid size difference compensation,
snapshots, SLC

plasmid and iMGE (multiple) FASTA-file export

+ rep-, iMGE-, and AMR-target(s) information

Plasmid Visualization

**Comparative Plasmid
Visualization**

pyGenomeViz

bioRxiv preprint doi: <https://doi.org/10.1101/2024.07.09.602722>; this version posted August 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

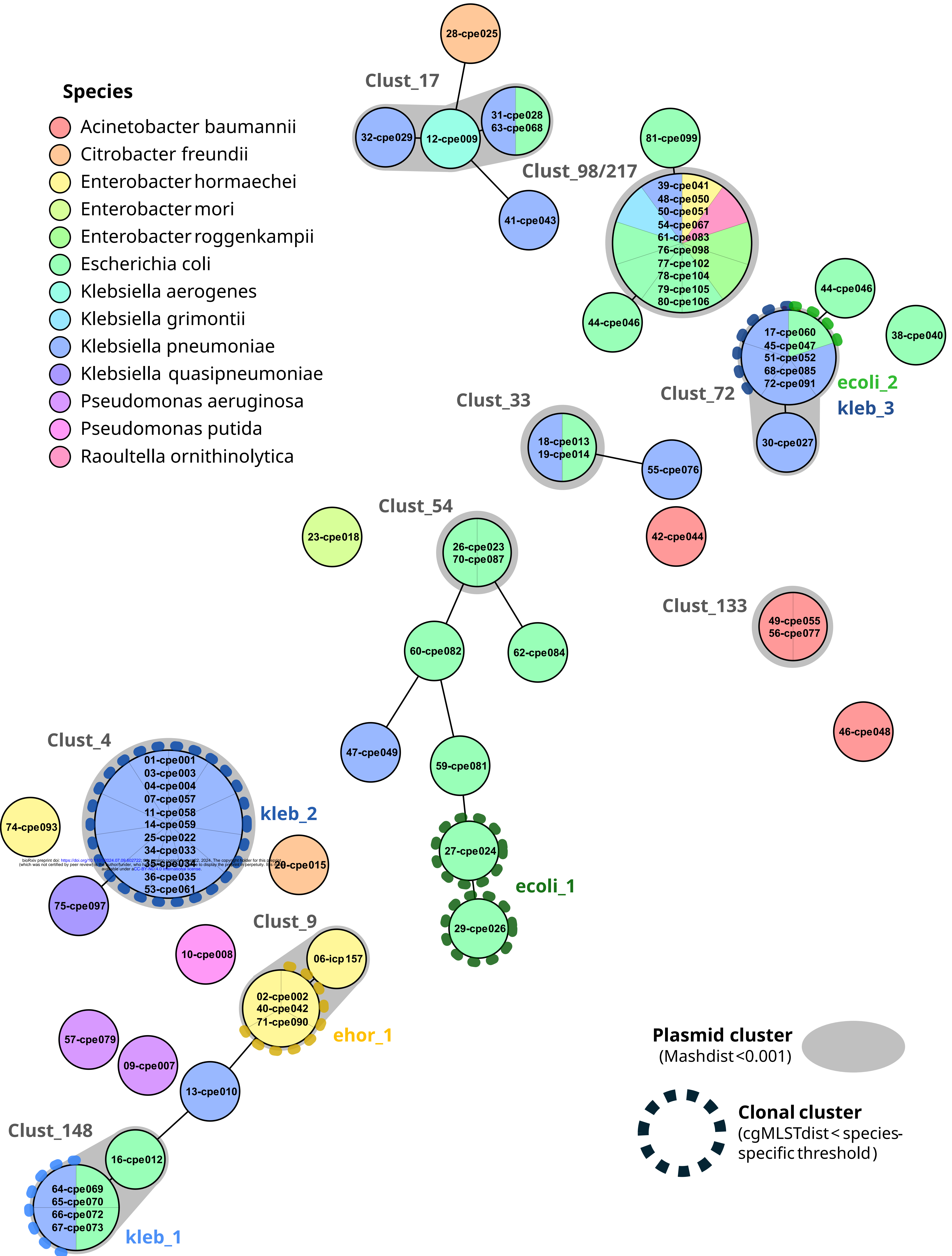
*Long-read data

**Single Linkage Clustering

*** Early Warning Alert

Species

- Acinetobacter baumannii
- Citrobacter freundii
- Enterobacter hormaechei
- Enterobacter mori
- Enterobacter roggenkampii
- Escherichia coli
- Klebsiella aerogenes
- Klebsiella grimontii
- Klebsiella pneumoniae
- Klebsiella quasipneumoniae
- Pseudomonas aeruginosa
- Pseudomonas putida
- Raoultella ornithinolytica



bioRxiv preprint doi: <https://doi.org/10.1101/2024.07.09.602722>; this version posted July 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

bioRxiv preprint doi: <https://doi.org/10.1101/2024.07.09.602722>; this version posted August 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

