

1 **Linear plasmids in *Klebsiella* and other Enterobacteriaceae**

2
3 Jane Hawkey¹, Hugh Cottingham¹, Alex Tokolyi², Ryan R. Wick¹, Louise M. Judd¹, Louise
4 Cerdeira³, Doroti de Oliveira Garcia⁴, Kelly L. Wyres¹, Kathryn E. Holt^{1,5}

5
6 1 Department of Infectious Diseases, Central Clinical School, Monash University,
7 Melbourne, Victoria 3004, Australia

8 2 Department of Human Genetics, Wellcome Sanger Institute, Hinxton, UK

9 3 Liverpool School of Tropical Medicine, Liverpool L3 5QA, UK

10 4 Adolfo Lutz Institute, Regional Laboratory Center, Marilia, Brazil

11 5 London School of Hygiene & Tropical Medicine, London WC1E 7HT, UK

12 **Corresponding authors**

13 jane.hawkey@monash.edu

14 kat.holt@lshtm.ac.uk

15 **Keywords**

16
17 linear plasmids, *Klebsiella*, Enterobacteriaceae, toxin-antitoxin

18 **Abstract**

19
20
21 Linear plasmids are extrachromosomal DNA that have been found in a small number of
22 bacterial species. To date, the only linear plasmids described in the Enterobacteriaceae
23 family belong to *Salmonella*, first found in *Salmonella* Typhi. Here, we describe a collection
24 of 12 isolates of the *Klebsiella pneumoniae* species complex in which we identified linear
25 plasmids. We used this collection to search public sequence databases and discovered an
26 additional 74 linear plasmid sequences in a variety of Enterobacteriaceae species. Gene
27 content analysis divided these plasmids into five distinct phylogroups, with very few genes
28 shared across more than two phylogroups. The majority of linear plasmid-encoded genes
29 are of unknown function, however each phylogroup carried its own unique toxin-antitoxin
30 system and genes with homology to those encoding the ParAB plasmid stability system.
31 Passage *in vitro* of the 12 linear plasmid-carrying *Klebsiella* isolates in our collection (which
32 include representatives of all five phylogroups) indicated that these linear plasmids can be
33 stably maintained, and our data suggest they can transmit between *K. pneumoniae* strains
34 (including members of globally disseminated multidrug resistant clones) and also between
35 diverse Enterobacteriaceae species. The linear plasmid sequences, and representative
36 isolates harbouring them, are made available as a resource to facilitate future studies on the
37 evolution and function of these novel plasmids.

38
39

40 **Significance as a BioResource to the community**

41 This study provides the first report of linear plasmids identified within the *Klebsiella*
42 *pneumoniae* species complex and the first report in Enterobacteriaceae besides *Salmonella*.
43 We present the first comparative analysis of linear plasmid sequences in
44 Enterobacteriaceae, however whilst this family is highly clinically significant, the functional
45 and/or evolutionary importance of these plasmids is not yet clear. To facilitate future studies
46 to address these questions, we have publicly deposited (i) the collection of linear plasmid
47 sequence data; (ii) isolates representative of each of the distinct linear plasmid phylogroups.
48

49 **Data Summary**

50 **The authors confirm all supporting data, code and protocols have been provided**
51 **within the article or through supplementary information.**

- 52 1. Whole genome sequence reads from *Klebsiella pneumoniae* isolates sequenced in
53 this study have been deposited in NCBI SRA under the accession numbers listed in
54 **Table S1**.
- 55 2. Representative annotated sequences of one linear plasmid per phylogroup have
56 been deposited in FigShare, **doi 10.26180/16729126**.
- 57 3. A copy of all linear plasmid sequences that we assembled from publicly available
58 genome sequence reads are available in FigShare, **doi 10.26180/16531365**. Read
59 accessions for these are given in **Table S1**.
- 60 4. Eleven representative *K. pneumoniae* isolates harbouring linear plasmids described
61 in this study have been deposited with the National Collection of Type Cultures
62 (NCTC) and are available for purchase under the NCTC accession numbers listed in
63 **Table S1**. *K. pneumoniae* 1194/11 (representative of phylogroup B) has been
64 deposited in the Microorganisms Collection Center, Adolfo Lutz Institute, São Paulo,
65 Brazil. To request strain 1194/11 (IAL 3063, SISGEN ABBF09B), contact:
66 Microorganisms Collection Center
67 Culture Collection Laboratory
68 Instituto Adolfo Lutz, Sao Paulo State Department of Health
69 Address: Av Dr Arnaldo, 351, 10 floor, room 1020
70 Phone number: +55 11 3068-2884
71 Zip code 01246-000, São Paulo, Brazil
72 E-mail: coleccion@ial.sp.gov.br
- 73 5. Alignments of terminal inverted repeat sequences for each phylogroup can be found
74 in **Data S1**, available on FigShare, **doi 10.26180/16531371**.

75

76 Introduction

77 Plasmids are extrachromosomal DNA that are frequently found in bacterial cells. The vast
78 majority of plasmid molecules exist in a circular conformation, however linear plasmids have
79 been found in several bacterial species, with the first description in *Streptomyces* in 1979 [1],
80 and later in *Borrelia* [2], where they are universally present. A study of clinical *Enterococcus*
81 *faecium* isolates recently reported the existence of a 143 kbp linear plasmid carrying a N-
82 acetyl-galactosamine (GalNAc) utilization operon that could be transferred between strains
83 via conjugation [3]. Linear plasmids appear to be exceedingly rare within
84 Enterobacteriaceae, with the first, pBSSB1 (27 kbp), described in 2007 from *Salmonella*
85 Typhi isolated in Indonesia [4]. Prior to this discovery, the only other linear replicons
86 described within Enterobacteriaceae were those derived from bacteriophage, including
87 pKO2 in *Klebsiella oxytoca* [5], N15 in *Escherichia coli* [6], and PY54 in *Yersinia*
88 *enterocolitica* [7]. These bacteriophage-derived linear replicons are distinct from the true
89 linear plasmids described in *Salmonella*, *Enterococcus*, *Streptomyces* and *Borrelia*, as they
90 still possess bacteriophage-specific genes including those for the lysis pathways [5].

91
92 For replicons that are linear, there is a requirement to stabilise the terminal ends to ensure
93 stability and appropriate replication, which in eukaryotes is achieved through the use of
94 telomeres. In contrast, bacterial linear plasmids can either (i) create hairpin loops, as in
95 *Borrelia* [8] and *Enterococcus* [3], or (ii) bind telomere-associated proteins to each end of the
96 molecule with the assistance of terminal inverted repeats (TIRs), as in *Streptomyces* [9]. The
97 *Salmonella* linear plasmid pBSSB1 was found to carry 1,230 bp TIRs with covalently bound
98 proteins on the end, similar to *Streptomyces*, however these had no homology to any
99 previously identified TIRs [4].

100
101 The *S. Typhi* linear plasmid pBSSB1 encodes two flagellar genes, an *fljA*-like gene and
102 *fljB^{z66}* [4]. *FljB^{z66}* encodes the phase II z66 flagellin antigen, whilst the *fljA*-like gene is
103 thought to encode the repressor of the chromosomally-encoded phase I flagellin antigen,
104 allowing for phase II z66 antigen presentation [4]. Few other genes from the 27 kbp plasmid
105 pBSSB1 have been characterised. Linear plasmids homologous to pBSSB1 have since been
106 described in other *Salmonella* serovars, at a prevalence of ~0.3%, the majority of which
107 carried the z66 flagellin genes [10].

108
109 In this study we report the discovery of multiple diverse linear plasmids in genomes
110 belonging to the *Klebsiella pneumoniae* species complex (*K. pneumoniae* and six closely
111 related taxa) within the Enterobacteriaceae. We demonstrate the linearity of these replicons
112 using long-read and short-read sequencing, show they are reliably maintained within their
113 natural host isolates during 10 rounds of laboratory passage, and identify homologs in the
114 genomes of several other Enterobacteriaceae species. Clustering on the basis of gene
115 content, we identify five major phylogroups of *K. pneumoniae* linear plasmids and describe
116 their sequence characteristics in terms of size, GC content, TIR sequence and TIR length.

117

118

119 **Methods**

120 ***Identifying linear plasmids in K. pneumoniae species complex genomes***

121 We screened for linear plasmids in the assembly graphs of 1,119 genomes of the *K.*
122 *pneumoniae* species complex, including 452 from our own collection of human clinical and
123 carriage isolates [11–13] and 667 publicly available read sets (see **Table 1**). Paired-end
124 Illumina reads for each genome were assembled using Unicycler v0.4.7, using default
125 parameters. The first assembly graph produced by Unicycler (001_best_spades_graph.gfa)
126 was searched for the signature two-contig structure of a linear plasmid (a connected
127 component of the graph consisting of one contig connected at both ends to the same end of
128 another contig, see **Figure 1a**) using a custom Python script (available at doi
129 10.26180/16531374). We subsequently used these linear plasmid sequences as queries for
130 a nucleotide BLAST search of the 1,119 genome assemblies, to recover instances where the
131 linear plasmid sequence was present but had not fully assembled into the characteristic two-
132 contig graph structure. This resulted in a total of 25 linear plasmid sequences, these have
133 been deposited in FigShare, doi 10.26180/16531365.

134

135 ***Identifying homologs in other species***

136 To detect homologous linear plasmids in other bacterial species, we performed a nucleotide
137 BLAST search of NCBI (May 10th, 2021), using as queries each of the linear plasmid
138 sequences identified in *Klebsiella*, as well as the pBSSB1 sequence (accession
139 NC_011422). Hits with $\geq 90\%$ identity and $\geq 60\%$ coverage of a query sequence were
140 considered as putative linear plasmid sequences (n=61). Metadata for each linear plasmid
141 sequence and its host bacterium was pulled from the GenBank record for the corresponding
142 whole genome sequence (WGS). To confirm the taxonomy and multi-locus sequence types
143 (MLST) of the bacterial hosts of these putative linear plasmids, the chromosome sequence
144 for each genome was uploaded to Pathogenwatch (<https://pathogen.watch>). For strain WP3-
145 W18-ESBL-02 (in which plasmid 3, accession AP021975.1, was a hit to linear plasmid query
146 pINF007 plasmid 3), Pathogenwatch was unable to detect a species, however the Genome-
147 based Taxonomy Database (using GTDB-Tk [14] with database release 202,
148 <https://gtdb.ecogenomic.org>) assigned it as a novel *Kluyvera* species, *Kluyvera ascorbata_B*.
149 **Table S1** lists the species given by the submitter in GenBank, in addition to species detected
150 by Pathogenwatch or GTDB, for all genomes.

151

152 ***Plasmid stability analysis***

153 For the 12 bacterial isolates in our collection with linear plasmids, we tested the stability of
154 these plasmids during 10 passages in broth culture. Isolates from frozen glycerol stocks
155 were streaked onto cation adjusted Mueller Hinton (CAMH) agar plates and incubated for 20
156 hours at 37°C. A single colony from each plate was streaked onto a fresh CAMH plate and
157 inoculated into 3 mL of CAMH broth, and both were incubated for 20 hours at 37°C. From
158 the broth culture, a glycerol stock and bacterial pellet, day 1 (D1) samples were prepared.
159 This process was repeated 9 additional times to yield day 2-10 (D2-D10) samples.

160

161 Long-read sequencing (Oxford Nanopore Technologies, ONT) was performed as previously
162 described [15]. Briefly, genomic DNA was prepared from the D1 and D10 bacterial pellets
163 using GenFind v3 reagents (Beckman Coulter). A long-read sequencing library was
164 prepared using the ligation library kit (LSK-109, ONT) with native barcoding expansion pack
165 (EXP-NBD104 and NBD114, ONT). The library was run on a R9.4.1 MinION flow cell for 48

166 hours yielding 2.75 Gbp of reads. Reads were base called with Guppy v3.3.3 using the
167 dna_r9.4.1_450bps_hac (high-accuracy) basecalling model.

168

169 To determine presence/absence and copy number of all plasmids in each genome, reads
170 were mapped to their respective reference genome assemblies (listed in **Table S1**) using
171 minimap2 v2.17 [16]. Mean read depth across each replicon in the assembly was calculated
172 using the read alignments, and copy number for each plasmid was determined by dividing
173 mean read depth across a plasmid replicon by the mean read depth across the
174 chromosome.

175

176 ***Confirming plasmid linearity***

177 For the 12 linear-plasmid-positive isolates in our collection, reads from the day 1 (D1) ONT
178 sequencing (see above) were aligned to their respective reference genomes using minimap2
179 v2.17 [16]. For each linear and circular plasmid sequence, we extracted all high quality read
180 alignments (read identity $\geq 80\%$, alignment length $\geq 1,000$ bp) that aligned within 90 bp of the
181 end of the plasmid reference sequence. For these reads, we calculated the proportion that
182 extended ≥ 100 bp beyond the edge of the plasmid reference sequence (and thus were soft-
183 clipped ≥ 100 bp by the read aligner). If the replicon from which the reads originated was
184 linear, we would expect to see few or no such soft-clipped reads, because the N'- and C'-
185 terminal ends of the plasmid ssDNA molecules should match the start and end of the
186 reference sequence (see **Figure 1b**). However, if the plasmid from which the reads
187 originated was circular, we would expect to see many reads that are soft-clipped at the ends
188 of the linearised reference sequence (see **Figure 1c**).

189

190 ***Linear plasmid characteristics and relationships***

191 To compare gene content across plasmid sequences, all 86 linear plasmid sequences
192 retrieved from Enterobacteriaceae genomes were annotated using Prokka v13.3 [17], and
193 genes were clustered into homologous groups using panaroo v1.2.4 [18], with a threshold of
194 70% amino acid identity to determine homology (details of clusters can be found in **Table**
195 **S2**). The panaroo gene presence/absence matrix (**Table S3**) was subjected to hierarchical
196 clustering using *hclust* in R (with default settings, i.e. Euclidean distance and *ward.D2*
197 clustering algorithm) to generate a dendrogram, which was cut into five phylogroups after
198 visual inspection.

199

200 TIR length was calculated by taking each linear plasmid sequence, obtaining the reverse
201 complement, and determining the length of sequence from the start of the forward and
202 reverse complement sequences that were identical, with zero mismatches. Nearly all (except
203 five) linear plasmid assemblies identified via nucleotide BLAST search of NCBI had very
204 small TIRs using this method ($n=56$, between 0-54 bp). We assume this is the result of
205 artefacts in the assembly process, which we are unable to explore without the underlying
206 sequence reads; therefore plasmid sequences available only as publicly deposited
207 assemblies without short reads were excluded from TIR length analyses. TIR sequences for
208 the 25 linear plasmids generated from our assemblies were extracted, categorised into their
209 respective phylogroups, and aligned using the clustalo algorithm in SeaView [19] to identify
210 regions of homology within phylogroups (**Data S1**).

211

212 Nucleotide divergence between linear plasmid sequences was calculated by performing
213 pairwise BLASTn alignments between all pairs of plasmids in the same phylogroup, and
214 extracting the percent identity of the longest hit.

215

216 **Detailed annotation of representative linear plasmids**

217 To further explore gene function in these linear plasmids, we undertook detailed annotation
218 for one representative per phylogroup (A, INF019; B, 1194/11; C, INF102; D, INF007; E,
219 INF352). Each representative was annotated using the RASTtk pipeline [20–22]. We
220 screened for PFAM domains for all genes identified by RAST with hmmscan [23] via the
221 EMBL-EBI server using default parameters. Resulting Pfam domains for genes with hits are
222 listed in **Table S4**. To determine if any genes in the representative plasmids had homology
223 to genes found in the Enterobacteriaceae, protein sequences were extracted from the RAST
224 annotations and screened using BLASTp to the refseq_select database on NCBI, restricting
225 results to Enterobacteriaceae. Genes with at least 50% protein identity to those in the
226 Enterobacteriaceae were considered sufficiently similar to have a similar function.
227 Representative plasmid annotations have been deposited in GenBank, accessions can be
228 found in **Table S1**. To determine conservation of genes amongst plasmids in the same
229 phylogroup, RAST annotations were matched with the Prokka annotations from the panaroo
230 analysis.

231

232 **Trinucleotide profiles of linear plasmids and bacterial chromosomes**

233 To investigate the potential donors of the linear plasmids into Enterobacteriaceae, we used
234 *compseq* from the EMBOSS package [24] to calculate the frequencies of all possible
235 trinucleotides in each of our 12 *Klebsiella* linear plasmids, their host chromosomes, as well
236 as one representative per bacterial species (n=47,893) as defined by the GTDB database
237 release 202 [25, 26]. We created a distance matrix using these frequencies with the *rdist*
238 function in the R package *fields* (<https://github.com/NCAR/Fields>).

239

240 **Results and Discussion**

241 **Identification of linear plasmids**

242 We identified unusual structures in the assembly graphs of some *K. pneumoniae* in our in-
243 house collection of genomes, which were consistent with linear plasmids with inverted
244 repeats at either end (**Figure 1a, Methods**). We systematically screened for these structures
245 in the assembly graphs of our in-house collection of *K. pneumoniae* species complex
246 isolates, collected from human clinical infections or colonisation [11, 12] in an Australian
247 hospital (n=452), as described in **Methods**. This screen yielded 11 genomes harbouring
248 linear plasmids (2.4% of genomes) including seven *K. pneumoniae* and four *Klebsiella*
249 *variicola* (**Table S1**). The corresponding isolates originated from nine patients, representing
250 three instances of asymptomatic colonisation (*K. pneumoniae* ST359, *K. variicola* ST386
251 and ST642), one instance of simultaneous gut colonisation and pneumonia (*K. pneumoniae*
252 ST37), and five instances of clinical infection (urinary tract infection with *K. pneumoniae*
253 ST20, ST27, ST1449; wound infection with *K. pneumoniae* ST3073 and *K. variicola* ST347).
254 The only extended-spectrum beta lactamase (ESBL)-positive (which confers resistance to
255 the third generation cephalosporins) isolates amongst those with identified linear plasmids
256 were two *K. variicola* ST347 isolated from the same patient nine days apart.

257

258 The linear plasmids were median 33,775 bp in size (range 31,739 - 44,271 bp), including the
259 TIRs at either end. To confirm our hypothesis that these plasmids were indeed linear
260 molecules, rather than the typical circular plasmid structure, we undertook additional
261 sequencing using long reads, and aligned the long reads to each linear plasmid (see
262 **Methods**). Plasmids were considered linear if there were few soft-clipped bases from reads
263 aligned at the start or end of the linear reference sequence (unlike a circular replicon, where
264 many reads are expected to overlap the ends of the linearised reference sequence, see
265 **Figure 1b**). The 12 linear plasmids had a median of 3.5% (range 0.6-32.4%) soft-clipped
266 start or end reads, compared to 98.5% (range 92.3-100%) for the circular plasmids (**Figure**
267 **1b, Fig S1, Table S1**). Additionally, all but two linear plasmids (those from *K. variicola*
268 ST347) were supported by reads (median n=70, range n=10 to 177) that spanned the full
269 length of the plasmid, including both TIRs (**Table S1**). Importantly, the soft-clipped parts of
270 the reads did not map to the other end of the plasmid sequence (as would be expected for a
271 circular plasmid), rather, they were chimeric reads, where two unrelated DNA segments
272 have fused during library preparation [27].

273

274 To investigate whether other linear plasmids are present in the *K. pneumoniae* species
275 complex, we generated and screened assembly graphs for an additional 667 publicly
276 available read sets, which represent a diverse set of (mostly human clinical) isolates from
277 multiple continents including Africa, Asia and Europe (**Table 1**). Across this set of genomes,
278 we identified linear plasmid graph structures in an additional 14 genomes (2.1%, see **Table**
279 **1**). The corresponding isolates include 12 *K. pneumoniae* from humans (UK, Kenya,
280 Cambodia, Brazil), one *K. pneumoniae* isolated from retail pork (USA), and one *Klebsiella*
281 *africana* human blood isolate (Kenya).

282

283 Using as queries the sequences of the 25 linear plasmids that we identified from *Klebsiella*
284 assembly graphs, we performed a BLAST search of the NCBI database to identify homologs
285 in other genomes (see **Methods**). This revealed another 61 putative linear plasmid
286 sequences; all were from Enterobacteriaceae, including *Klebsiella* (n=23, including 17 *K.*
287 *pneumoniae*), *Salmonella enterica* (n=16, including pBSSB1), *Citrobacter* (n=8),
288 *Enterobacter* (n=7), *Escherichia coli* (n=3), *Serratia marcescens* (n=2), *Phytobacter*
289 *diazotrophicus* (n=1) and *Kluyvera ascorbata_B* (n=1) (**Table S1**). Genomes harbouring
290 linear plasmids came from a wide variety of sources, including bacteria isolated from water
291 (n=19), humans (n=13), food (n=4), animals (n=3), and plants (n=1) (**Table S1**). Amongst
292 the linear-plasmid-positive *K. pneumoniae* were well-known carbapenemase-producing and
293 ESBL producing clones: ST340 (n=3, KPC-4 and CTX-M-15), ST258 (KPC-2 and SHV-12),
294 ST11 (n=1, KPC-2 and SHV-12), ST147 (n=1, OXA-181 and CTX-M-15). Hundreds of
295 genomes of each of these clones are present in the NCBI database and the vast majority do
296 not harbour linear plasmid sequences, suggesting that the linear-plasmid-positive variants
297 are rare, and likely result from recent horizontal transfer but this has not resulted in clonal
298 expansion during which the plasmid has been stably maintained. This is in contrast to the
299 recent report in *E. faecium* where the linear plasmid *pELF_USZ* was stably maintained in a
300 host lineage during >2 years of clonal spread in a hospital [3].

301

302 **Characteristics of linear plasmids in Enterobacteriaceae**

303 We compiled the full set of 86 linear plasmid sequences (25 identified from assembly
304 graphs, plus 61 inferred from homology via BLAST) and clustered them by their gene
305 content (see Methods). This revealed five distinct linear plasmid phylogroups (which we

306 labelled A-E, see **Figure 2, Table S2 & S3**), with very little gene sharing between
307 phylogroups (genes defined as homologous if they had >70% nucleotide identity). Each
308 phylogroup included sequences from multiple genera, notably all five phylogroups were
309 detected in both *Klebsiella* and *Salmonella* (**Figure 2**). No genes were present across more
310 than two phylogroups, but each phylogroup had a core set of genes found in $\geq 95\%$ of
311 plasmids in that group; these represented between 15% (phylogroup E) to 47% (phylogroup
312 B) of all genes found in that phylogroup (**Figure 3a, Figure 4a**). Nucleotide diversity within
313 phylogroups varied (**Figure 3b**), with phylogroup A displaying significantly greater pairwise
314 divergence across the full plasmid sequence than phylogroups B, C and D (mean 6%
315 divergence vs mean 2.6-4.2%, $p < 1 \times 10^{-16}$ using Wilcoxon test for A vs B, C or D).
316 Phylogroup E showed a high range in divergence (0-16%, mean 4.2%), due to the presence
317 of two divergent subgroups (see **Figure 2**).

318
319 The vast majority of genes annotated in each linear plasmid were hypothetical proteins and
320 had no close homologs in other Enterobacteriaceae genomes (**Figure 4b**). However, there
321 were few reference plasmid genes ($n=55$, 20%) for which we were able to obtain some form
322 of functional annotation based on sequence homology or protein domain matches (see
323 **Methods, Table S4**). Most of these annotations were for genes encoding proteins likely
324 relevant to basic plasmid maintenance functions. All five phylogroups carried genes with
325 type II toxin-antitoxin domains (see **Table S4**), which are often found on plasmids and can
326 enable plasmid maintenance by performing post segregational killing of daughter cells that
327 do not carry the plasmid [28]. These systems were core in all phylogroups. Phylogroups A,
328 B, C and D each carried a *relBE* family system (65%-83% homology between variants in
329 phylogroups B-D, A carried a distinct variant), whilst phylogroup E carried a *vapBC* system
330 (**Figure 4b**). These toxin-antitoxin clusters generally had at least one gene of the pair
331 encoding a protein with $\geq 50\%$ homology to toxin-antitoxin systems found in
332 Enterobacteriaceae (**Figure 4b, Table S4**). Pairs of adjacent genes encoding novel proteins
333 with Pfam matches to the partitioning proteins ParA (PF13614 or PF01656) and ParB
334 (PF18821) were detected as core in each phylogroup (**Figure 4, Table S4**). These likely
335 contribute to control of plasmid segregation into daughter cells [29], homologous sequences
336 were not detected in other Enterobacteriaceae. Sequences with homology to the
337 transcriptional repressor *hns* were identified in all phylogroups except A (**Figure 4b**),
338 however the encoded proteins were highly divergent from one another (27-66%) and the
339 genes were classed as separate gene groups by panaroo (**Table S2**). *Hns* are commonly
340 plasmid-encoded and regulate expression of both plasmid and chromosomally-encoded
341 genes. Proteins with hits to known restriction/modification domains were also identified in all
342 reference plasmids, these are frequently encoded by mobile elements and can function as
343 toxin/antitoxin systems to force maintenance of those elements. Phylogroup A was the only
344 phylogroup in which flagellin genes were identified, in $n=7/16$ plasmid sequences. One of
345 these was plasmid pBSSB1, and the other six were all linear plasmids from *Salmonella*
346 *enterica* serovar Senftenberg isolated from Switzerland [10]. Phylogroups C and D both
347 carried three core genes apiece harbouring PilS (type IV pilin) domains (**Figure 4b**), which
348 could potentially function as adhesins.

349
350 All five phylogroups differed substantially from one another in their basic characteristics,
351 including plasmid length, TIR length and GC content. Phylogroups D and C had the longest
352 plasmids (medians 40.9 kbp and 42.9 kbp respectively), and phylogroup B the smallest
353 (median 23.7 kbp, **Figure 5a**). We calculated the size of TIRs by aligning the beginning of

354 each plasmid to the reverse complement of itself (see **Methods**). We were able to detect a
355 TIR in 57 of the linear plasmid sequences. Those without a TIR were all identified in publicly
356 available assemblies that were assembled using a variety of methods, and we hypothesise
357 that the lack of TIR sequence is most likely due to incomplete or fragmented assembly of the
358 plasmid, rather than lack of TIR in the sequenced molecules. For plasmids where we
359 performed the assembly in-house, we found that the length of the TIR differed substantially
360 between phylogroups, with phylogroups A and D having the longest TIRs (medians 1168 bp
361 and 1074 bp respectively), whilst phylogroups B, C and E had TIRs of approximately half
362 that length (medians 542 bp, 530 bp and 670 bp respectively, **Figure 5b**). There was a high
363 level of sequence conservation for TIRs within phylogroups, with a median of 89-97%
364 similarity in this region in phylogroups A - D (**Data S1, Figure S2**). Phylogroup E was more
365 diverse with an overall similarity of 65%; however, inspection of the alignments revealed that
366 this phylogroup carried two distinct TIR sequences, with a median of 89-99% TIR sequence
367 identity within each TIR grouping (**Data S1, Figure S2**). Finally, %GC for the linear plasmids
368 was very low in comparison to the normal chromosomal %GC range for Enterobacteriaceae,
369 which is typically ~50% (median 57% for the *Klebsiella* carrying linear plasmids). All linear
370 plasmid phylogroups had %GC <40%, with phylogroup B having the lowest out of all the
371 phylogroups (median 28%, compared to 34-35% for other phylogroups, $p < 2.5 \times 10^{-4}$ for all
372 comparisons, Wilcoxon test, **Figure 5c**).

373

374 **Potential donors of linear plasmids and their stability in *Klebsiella***

375 Given that linear plasmids are rare in *Klebsiella* and have a significantly lower %GC than
376 their host chromosomes, we assume that Enterobacteriaceae are unlikely to be the typical
377 hosts for these plasmids. We used trinucleotide frequencies as a genomic signature to
378 attempt to identify potential original hosts of these plasmids by calculating the distance
379 between our linear plasmids, their *Klebsiella* host chromosomes, and one representative per
380 bacterial species defined in the GTDB (see **Methods**). The 12 *Klebsiella* linear plasmids
381 clustered separately from their corresponding host chromosomes, with a mean distance of
382 2.3 between the chromosomes and linear plasmids (**Fig S3**). The *Klebsiella* chromosomes
383 were much more similar to each other than the linear plasmids (mean pairwise distance of
384 0.07 between chromosomal sequences vs 1.27 between pairs of linear plasmid sequences),
385 and clustered closely with other representatives of *Klebsiella* in the GTDB database (nearest
386 neighbour accession GCF_000742135.1, distance 0.09). The linear plasmid from 1194/11
387 clustered most closely to the Firmicute DUOC01 sp012839065 (accession
388 GCA_012839065.1, distance 1.3). This organism belongs to a strain from the class
389 Thermosediminibacteria, which was detected in a metagenomic sample obtained from an
390 anaerobic digester [30]. The other 11 linear plasmids were their own nearest neighbours
391 (median pairwise distance 1.09), the closest GTDB profile was Proteobacteria isolate
392 Neptuniibacter sp002435145 (accession GCA_002435145.1, mean distance 1.15 to the 11
393 plasmids). This organism belongs to the order Pseudomonadales, and was detected in a
394 marine environment [31].

395

396 To understand whether linear plasmids could be stably maintained within *Klebsiella*, we
397 undertook passage experiments on the 11 *Klebsiella* genomes carrying linear plasmids in
398 our collection. We performed long-read sequencing on all parental isolates (D1), passaged
399 each isolate 10 times (one passage per 24 hour period), and then performed long-read
400 sequencing on the final D10 isolates (see **Methods**). We found that all plasmids, both linear
401 and circular, were maintained in all genomes across 10 passages (**Figure 6**). Linear plasmid

402 copy number was generally estimated at ~1 per cell at both D1 and D10, with the exception
403 of 1194/11 (the only representative of phylogroup B), which had a copy number of 2-4, and
404 two of the phylogroup E plasmids (strains INF345, INF352) with copy number ~2 (see
405 **Figure 6**).

406

407 **Conclusions**

408 Here we provide the first (to our knowledge) collection of linear plasmids in the *K.*
409 *pneumoniae* species complex alongside a detailed description of their characteristics. Our
410 data show these plasmids are uncommon in *Klebsiella* and other *Enterobacteriaceae*
411 species, but can be stably maintained and can transfer between distinct *K. pneumoniae*
412 strains (including representatives of the globally-distributed multidrug resistant clones) and
413 other diverse *Enterobacteriaceae* [4]. The novel *Klebsiella* linear plasmids described here do
414 not carry any known antimicrobial resistance, virulence or metabolic genes; however
415 carriage of a linear plasmid has previously been shown to provide a metabolic advantage for
416 vancomycin-resistant *Enterococcus faecium* in the human gut [3] and to enable flagellar
417 antigen switching in *Salmonella* Typhi. By making freely available these linear plasmid
418 sequences and representative isolates that carry them, we hope to facilitate future research
419 into the function and potential evolutionary or clinical significance of these enigmatic
420 replicons.

421

422 **Authors and contributions**

423 Conceptualization, R.R.W, K.L.W and K.E.H; Formal analysis, J.H, H.C, A.T, R.R.W, L.M.J,
424 K.E.H; Methodology, R.R.W, A.T, L.M.J, K.L.W and K.E.H; Software, A.T and R.R.W;
425 Resources, L.C and D.O.G; Visualization, J.H, H.C, R.R.W and K.E.H; Writing - Original
426 Draft, J.H, H.C and K.E.H; Writing - Review & Editing, all authors; Funding acquisition,
427 K.E.H; Project administration, K.E.H.

428

429 **Conflicts of interest**

430 The authors declare that there are no conflicts of interest.

431

432 **Funding information**

433 K.E.H. was supported by a Senior Medical Research Fellowship from the Viertel Foundation
434 of Australia. This work was supported in part by the Bill & Melinda Gates Foundation
435 OPP1175797. Under the grant conditions of the Foundation, a Creative Commons
436 Attribution 4.0 Generic License has already been assigned to the Author Accepted
437 Manuscript version that might arise from this submission.

438

439

440 References

- 441 1. **Hayakawa T, Tanaka T, Sakaguchi K, Otake N, Yonehara H.** A linear plasmid-like DNA
442 in *Streptomyces* sp. producing lankacidin group antibiotics. *J Gen Appl Microbiol*
443 1979;25:255–260.
- 444 2. **Plasterk RHA, Simon MI, Barbour AG.** Transposition of structural genes to an
445 expression sequence on a linear plasmid causes antigenic variation in the bacterium *Borrelia*
446 *hermsii*. *Nature* 1985;318:257–263.
- 447 3. **Boumassoud M, Haunreiter VD, Schweizer TA, Meyer L, Chakrakodi B, et al.**
448 Genomic surveillance of vancomycin-resistant *Enterococcus faecium* reveals spread of a
449 linear plasmid conferring a nutrient utilization advantage. *bioRxiv* 2021;2021.05.07.442932.
- 450 4. **Baker S, Hardy J, Sanderson KE, Quail M, Goodhead I, et al.** A Novel Linear Plasmid
451 Mediates Flagellar Variation in *Salmonella* Typhi. *Plos Pathog* 2007;3:e59.
- 452 5. **Casjens SR, Gilcrease EB, Huang WM, Bunny KL, Pedulla ML, et al.** The pKO2 Linear
453 Plasmid Prophage of *Klebsiella oxytoca*. *J Bacteriol* 2004;186:1818–1832.
- 454 6. **Ravin V, Ravin N, Casjens S, Ford ME, Hatfull GF, et al.** Genomic sequence and
455 analysis of the atypical temperate bacteriophage N15. *J Mol Biol* 2000;299:53–73.
- 456 7. **Hertwig S, Klein I, Lurz R, Lanka E, Appel B.** PY54, a linear plasmid prophage of
457 *Yersinia enterocolitica* with covalently closed ends. *Mol Microbiol* 2003;48:989–1003.
- 458 8. **Lucyshyn D, Huang SH, Kobryn K.** Spring loading a pre-cleavage intermediate for
459 hairpin telomere formation. *Nucleic Acids Res* 2015;43:6062–6074.
- 460 9. **Yang C-C, Tseng S-M, Chen CW.** Telomere-associated proteins add deoxynucleotides to
461 terminal proteins during replication of the telomeres of linear chromosomes and plasmids in
462 *Streptomyces*. *Nucleic Acids Res* 2015;43:6373–6383.
- 463 10. **Robertson J, Lin J, Wren-Hedgus A, Arya G, Carrillo C, et al.** Development of a multi-
464 locus typing scheme for an Enterobacteriaceae linear plasmid that mediates inter-species
465 transfer of flagella. *Plos One* 2019;14:e0218638.
- 466 11. **Gorrie CL, Mirceta M, Wick RR, Edwards DJ, Strugnell RA, et al.** Gastrointestinal
467 carriage is a major reservoir of *K. pneumoniae* infection in intensive care patients. *bioRxiv*
468 2017;096446–096446.
- 469 12. **Gorrie CL, Mirceta M, Wick RR, Judd LM, Wyres KL, et al.** Antimicrobial-resistant
470 *Klebsiella pneumoniae* carriage and infection in specialized geriatric care wards linked to
471 acquisition in the referring hospital. *Clin Infect Dis* 2018;67:161–170.
- 472 13. **Bueno MFC, Francisco GR, O'Hara JA, Garcia D de O, Doi Y.** Coproduction of 16S
473 rRNA Methyltransferase RmtD or RmtG with KPC-2 and CTX-M Group Extended-Spectrum
474 β -Lactamases in *Klebsiella pneumoniae*. *Antimicrob Agents Ch* 2013;57:2397–2400.
- 475 14. **Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH.** GTDB-Tk: a toolkit to classify
476 genomes with the Genome Taxonomy Database. *Bioinform Oxf Engl*.
477 2019;10.1093/bioinformatics/btz848.
- 478 15. **Wick RR, Judd LM, Gorrie CL, Holt KE.** Completing bacterial genome assemblies with
479 multiplex MinION sequencing. *Microb Genom* 2017;3:e000132–e000132.
- 480 16. **Li H.** Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
481 2018;34:3094–3100.
- 482 17. **Seemann T.** Prokka: rapid prokaryotic genome annotation. *Bioinformatics*
483 2014;30:2068–2069.
- 484 18. **Tonkin-Hill G, MacAlasdair N, Ruis C, Weimann A, Horesh G, et al.** Producing
485 polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biol* 2020;21:180.
- 486 19. **Gouy M, Guindon S, Gascuel O.** SeaView version 4: A multiplatform graphical user

- 487 interface for sequence alignment and phylogenetic tree building. *Mol Biol Evol* 2009;27:221–
488 4.
- 489 20. **Aziz RK, Bartels D, Best AA, DeJongh M, Disz T**, et al. The RAST Server: Rapid
490 Annotations using Subsystems Technology. *BMC Genomics* 2008;9:75–75.
- 491 21. **Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ**, et al. The SEED and the Rapid
492 Annotation of microbial genomes using Subsystems Technology (RAST). *Nucleic Acids Res*
493 2014;42:D206 14-D206 14.
- 494 22. **Brettin T, Davis JJ, Disz T, Edwards RA, Gerdes S**, et al. RASTtk: A modular and
495 extensible implementation of the RAST algorithm for building custom annotation pipelines
496 and annotating batches of genomes. *Sci Rep* 2015;5:8365.
- 497 23. **Finn RD, Clements J, Eddy SR**. HMMER web server: interactive sequence similarity
498 searching. *Nucleic Acids Res* 2011;39:W29–W37.
- 499 24. **Rice P, Longden I, Bleasby A**. EMBOSS: the European Molecular Biology Open
500 Software Suite. *Trends Genetics* 2000;16:276–7.
- 501 25. **Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ**, et al. A complete
502 domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* 2020;38:1079–1086.
- 503 26. **Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A**, et al. A standardized
504 bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat*
505 *Biotechnol* 2018;36:996–1004.
- 506 27. **Wick RR, Judd LM, Wyres KL, Holt KE**. Recovery of small plasmid sequences via
507 Oxford Nanopore sequencing. *Microb Genom* 2021;7:10.1099/mgen.0.000631.
- 508 28. **Fraikin N, Goormaghtigh F, Melderer LV**. Type II Toxin-Antitoxin Systems: Evolution
509 and Revolutions. *J Bacteriol* 2020;202(7): e00763-19.
- 510 29. **Roberts MAJ, Wadhams GH, Hadfield KA, Tickner S, Armitage JP**. ParA-like protein
511 uses nonspecific chromosomal DNA binding to partition protein complexes. *Proc National*
512 *Acad Sci* 2012;109:6698–6703.
- 513 30. **Campanaro S, Treu L, Rodriguez-R LM, Kovalovszki A, Ziels RM**, et al. New insights
514 from the biogas microbiome by comprehensive genome-resolved metagenomics of nearly
515 1600 species originating from multiple anaerobic digesters. *Biotechnol Biofuels* 2020;13:25.
- 516 31. **Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ**, et al. Recovery of
517 nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat*
518 *Microbiol* 2017;2:1533–1542.
- 519 32. **Cerdeira L, Fernandes MR, Francisco GR, Bueno MFC, lenne S**, et al. Draft Genome
520 Sequence of a Hospital-Associated Clone of *Klebsiella pneumoniae* ST340/CC258
521 Coproducing RmtG and KPC-2 Isolated from a Pediatric Patient. *Genome Announc*
522 2016;4:e01130-16.
- 523 33. **Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH**, et al. Predicting antimicrobial
524 susceptibilities for *Escherichia coli* and *Klebsiella pneumoniae* isolates using whole genomic
525 sequence data. *J Antimicrob Chemoth* 2013;68:2234–2244.
- 526 34. **Smit PW, Stoesser N, Pol S, Kleef E van, Oonsivilai M**, et al. Transmission Dynamics
527 of Hyper-Endemic Multi-Drug Resistant *Klebsiella pneumoniae* in a Southeast Asian
528 Neonatal Unit: A Longitudinal Study With Whole Genome Sequencing. *Front Microbiol*
529 2018;9:1197.
- 530 35. **Davis GS, Waits K, Nordstrom L, Weaver B, Aziz M**, et al. Intermingled *Klebsiella*
531 *pneumoniae* Populations Between Retail Meats and Human Urinary Tract Infections. *Clin*
532 *Infect Dis* 2015;61:892–899.
- 533 36. **Henson SP, Boinett CJ, Ellington MJ, Kagia N, Mwarumba S**, et al. Molecular
534 epidemiology of *Klebsiella pneumoniae* invasive infections over a decade at Kilifi County

535 Hospital in Kenya. *Int J Med Microbiol* 2017;307:422–429.
536 37. **Moradigaravand D, Martin V, Peacock SJ, Parkhill J.** Evolution and Epidemiology of
537 Multidrug-Resistant *Klebsiella pneumoniae* in the United Kingdom and Ireland. *Mbio*
538 2017;8:e01976-16.

539 **Figures and Tables**

540

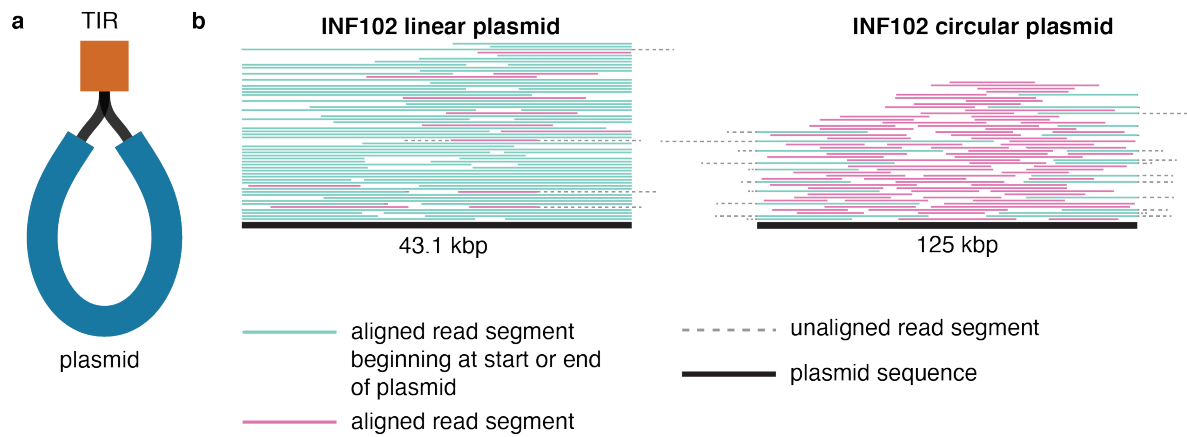
541 **Table 1:** Number of genomes positive for a linear plasmid across multiple different studies
542 from a variety of geographic regions and sampling types.

Dataset	# genomes	# linear plasmids	Country of origin	Sampling type
In-house collection (KASPAH) [11, 12]	452	11 (2.4%)	Australia	Humans (infections and faecal carriage)
Bueno 2013 [13, 32]	8	1 (12.5%)	Brazil	Humans, agricultural animals, urban waterways
Stoesser 2013 [33]	69	2 (2.9%)	UK	Humans (bloodstream infections)
Smit 2018 [34]	90	3 (3.3%)	Cambodia	Humans (neonatal care unit)
Davis 2015 [35]	61	1 (1.6%)	UK	Humans (urinary tract infections) and retail meat
Henson 2017 [36]	185	5 (2.7%)	Kenya	Humans (bloodstream infections)
Moradigaravand 2017 [37]	250	2 (0.8%)	UK and Ireland	Humans (bloodstream infections)

543

544

545



546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

Figure 1: Using read sequence data to determine linearity of plasmid sequences. a, Short read assembly graph structure of a linear plasmid. The plasmid consists of two contigs, the main plasmid sequence (blue, labelled 'plasmid'), connected to a second, shorter contig which is the terminal inverted repeat at both ends (orange, labelled 'TIR'). **b,** Long reads aligned to the linear and circular plasmid sequences from INF102, with the total number of alignments shown capped at 100 to improve visualisation. The plasmid sequence is the thick black line at the bottom, and reads aligning to the plasmid are shown in green if the alignment starts at the beginning or end of the plasmid sequence, or pink if the alignment starts elsewhere. Segments coloured dotted grey indicate regions of the read that do not align. Alignments to the linear plasmid have very few reads which soft-clip off the ends of the plasmid sequence, indicating linearity. Conversely, alignments to a circular plasmid have many reads soft-clipping over the edges of the plasmid sequence, indicating that this replicon is circular.

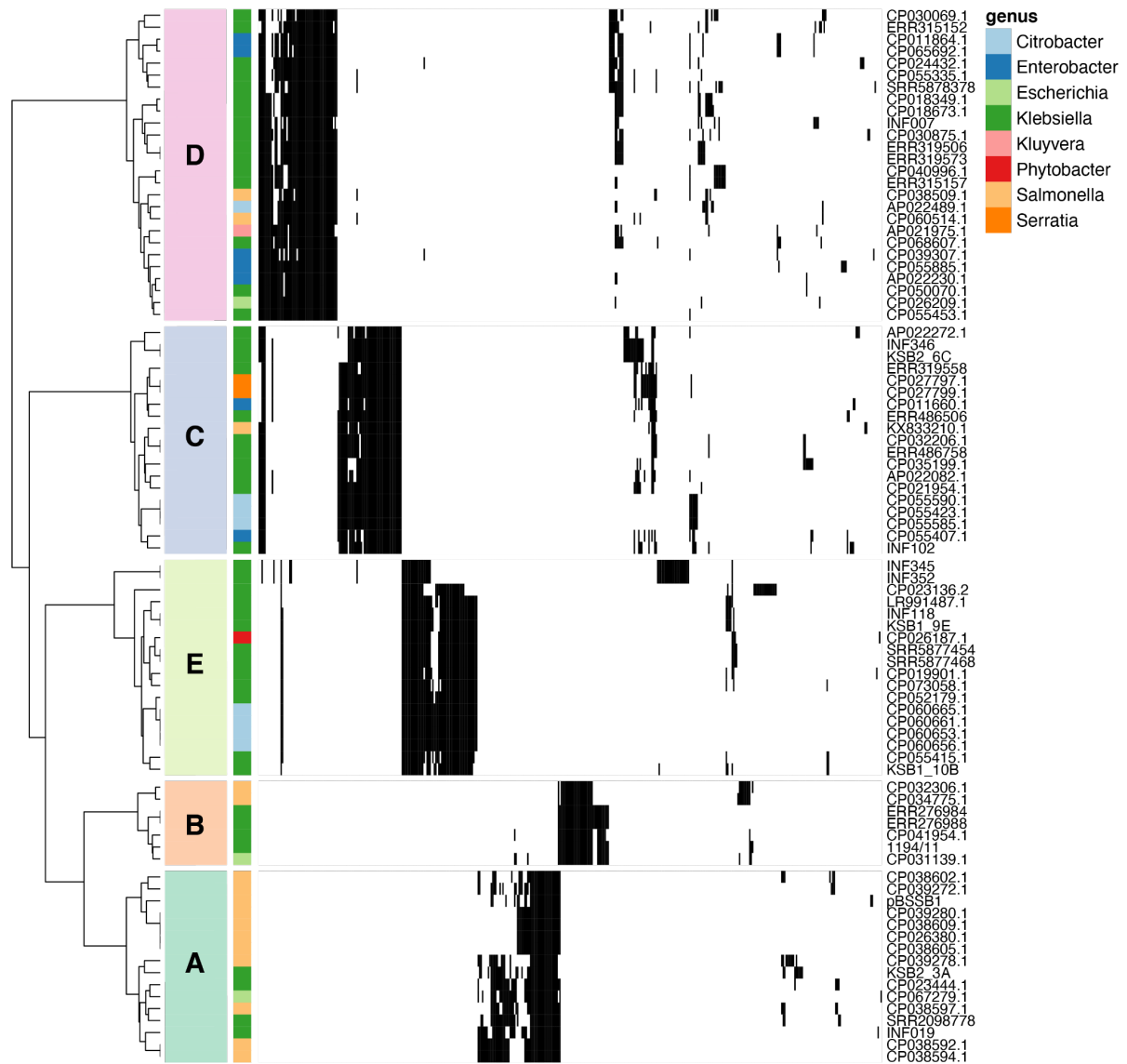
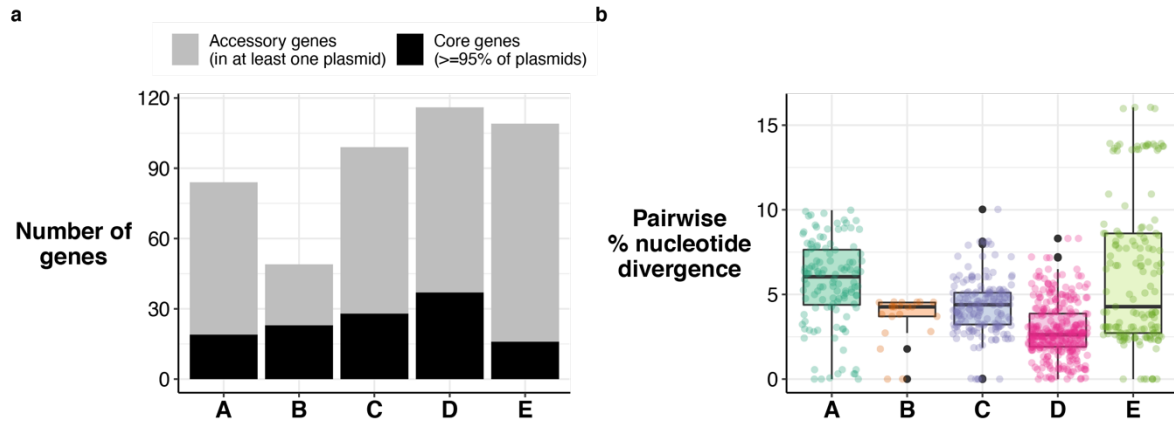


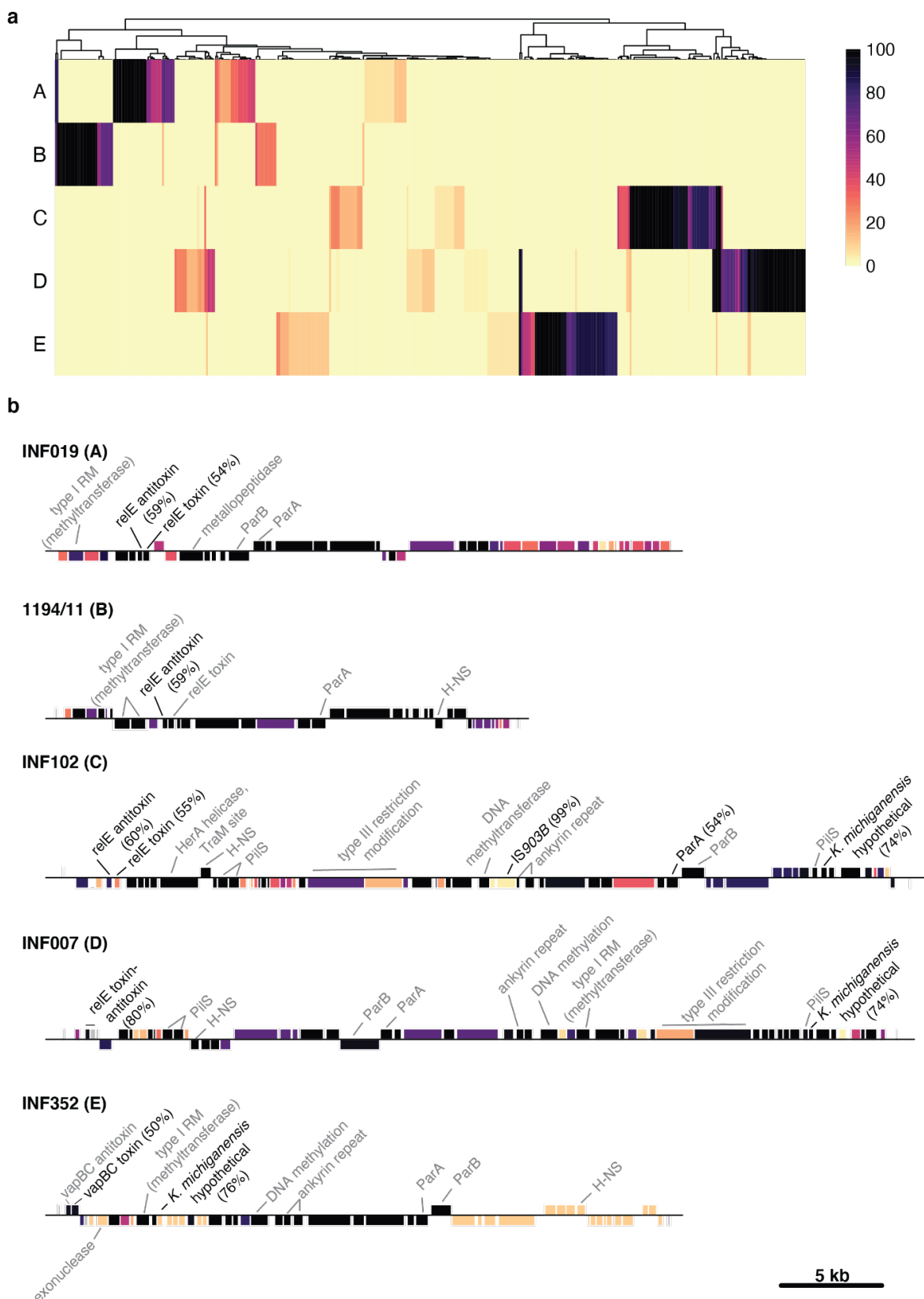
Figure 2: Hierarchical clustering of linear plasmids based on gene content. Plasmids were clustered with the *hclust* algorithm using the *ward.D2* method, and divided into five phylogroups (labelled in coloured boxes). Rows are annotated with the bacterial genus each linear plasmid was found in as per legend. Black indicates the presence of a gene, white absence. Plasmids are labelled with their names as per **Table S1**, and details of each gene can be found in **Table S2**.

568



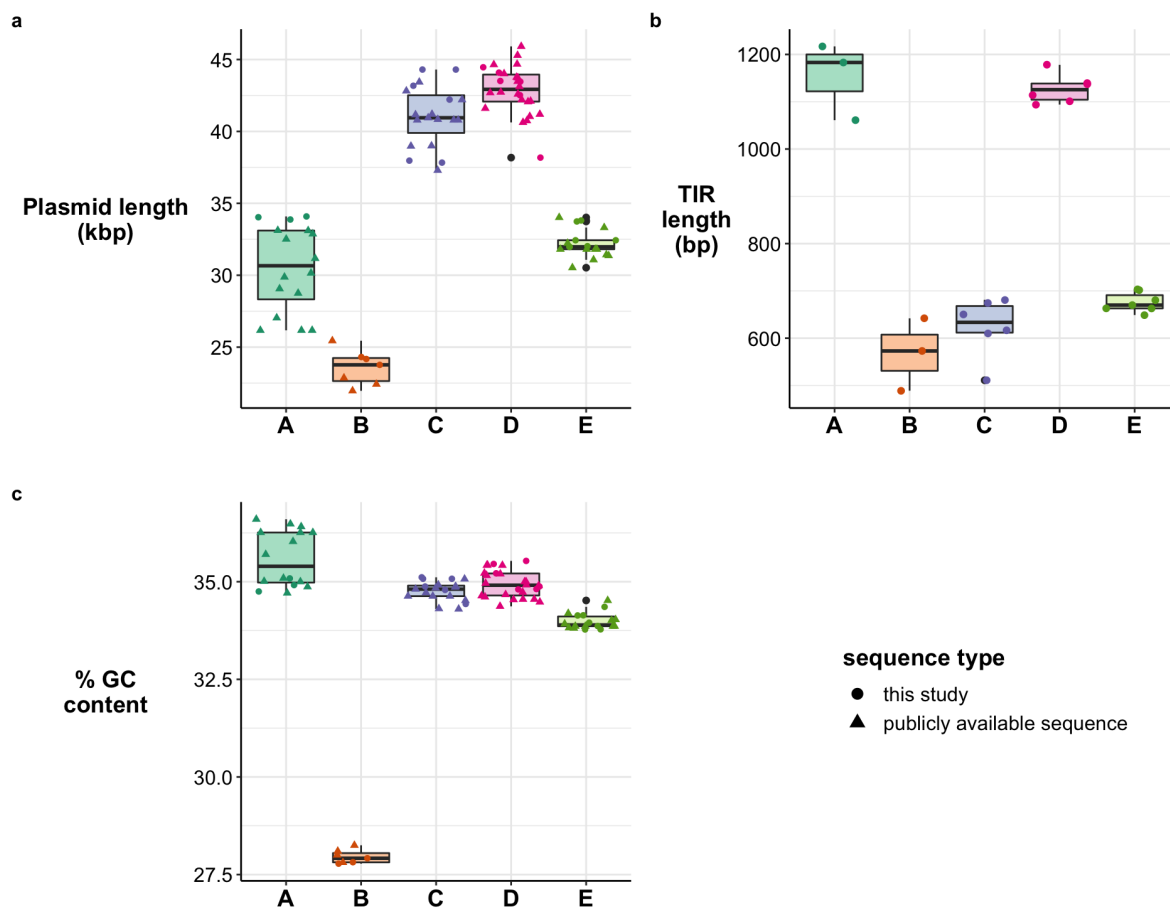
569
570
571
572
573
574
575
576
577
578

Figure 3: Core and accessory gene content by phylogroup, and nucleotide divergence by phylogroup. **a**, Number of core and accessory genes in each phylogroup. Bar height indicates the total number of genes found in at least one linear plasmid in each phylogroup. Black indicates the number of core genes (found in $\geq 95\%$ of plasmids); grey the number of accessory genes, as per legend. **b**, Distribution of pairwise nucleotide divergence within each phylogroup. Boxplots show median (thick black line), 1st and 3rd quartiles (edges of box), and solid lines give 1.5x the interquartile range. Outliers are shown as black dots. Individual values are shown as coloured dots.

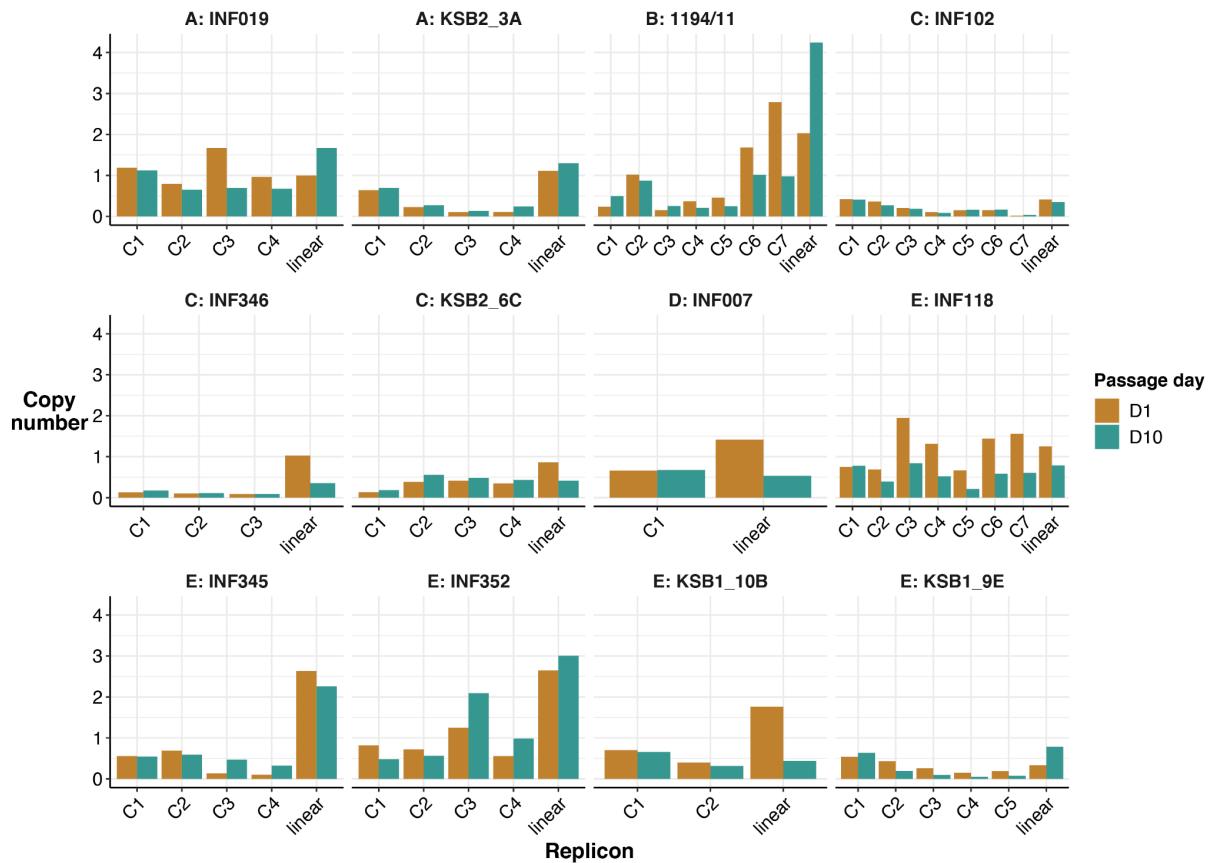


579
 580 **Figure 4: Conservation and function of genes in each phylogroup.** **a**, Heatmap showing
 581 proportion of plasmids with each gene by phylogroup. Columns are genes, clustered using
 582 *hclust*; rows are phylogroups (unclustered). Colour within each cell indicates the proportion
 583 of plasmids carrying each gene, as per legend. **b**, Gene maps of one representative plasmid

584 per phylogroup. Genes are indicated by blocks (above line - forward orientation, below line -
585 reverse orientation) and coloured by conservation in their phylogroup. Genes with $\geq 50\%$
586 homology to known genes in Enterobacteriaceae are indicated by black lines and text, with
587 gene homology shown in brackets. Genes with detected PFAM domains are indicated by
588 grey lines and text. Details of each gene can be found in **Table S4**.
589



590 **Figure 5: Characteristics of linear plasmid phylogroups. a, Distribution of plasmid**
591 **lengths in kbp.** Boxplots show median (black line), 1st and 3rd quartiles (edges of box), and
592 solid lines give 1.5x the interquartile range. Outliers are shown as black dots. Individual
593 values are shown as coloured dots or triangles, with shape indicating the origin of the
594 sequence as per legend. **b, Distribution of TIR lengths in bp,** as per (a). Publicly available
595 sequences are not represented in this plot due to assembly errors in the TIR region. **c,**
596 **Distribution of GC content,** as per (a).
597
598



599

600 **Figure 6: Estimated copy number of all plasmid replicons in each genome.** Height of
 601 each bar indicates copy number, coloured by passage day as per legend. Each pair of bars
 602 represents a plasmid, C[N] indicates a circular plasmid, Linear indicates the linear plasmid in
 603 that genome.

604

605 **Figure S1: Long read alignments to linear plasmids and one representative circular**
 606 **plasmid per *Klebsiella* genome.** The total number of alignments shown is capped at 100 to
 607 improve visualisation. The plasmid sequence is the thick black line at the bottom, and reads
 608 aligning to the plasmid are shown in green if the alignment starts at the beginning or end of
 609 the plasmid sequence, or pink if the alignment starts elsewhere. Segments coloured dotted
 610 grey indicate regions of the read that do not align.

611

612 **Figure S2: TIR sequence alignments within each phylogroup.** Linear plasmid sequences
 613 are clustered by gene content, with the phylogroup indicated by tip colour and coloured as
 614 per legend. TIR sequences are aligned within each phylogroup, where each colour
 615 represents a different nucleotide as per legend. Colour intensity indicates level of
 616 conservation at that position (pale=low; intense=high).

617

618 **Figure S3: Cluster dendrogram of trinucleotide frequencies for the 12 linear plasmids**
 619 **and their host *Klebsiella* chromosomes.** Trinucleotide frequencies were clustered using
 620 *hclust*. Tips are coloured by phylogroup or chromosome (as per legend).

621

622 **Table S1: Details of all linear plasmids described in this study.**

623

624 **Table S2: Details of panaroo pangenome analysis.**

625

626 **Table S3: Presence/absence of all genes in each linear plasmids.**

627

628 **Table S4: Gene annotation details for annotations in representative plasmids,**

629 **including PFAM hits and hits to *Enterobacteriaceae*.**