# A Peek into the Plasmidome of Global Sewage

**Philipp Kirstahler[1], Frederik Teudt[1], Saria Otani[1], Frank M. Aarestrup[1], and Sünje Johanna Pamp[1*]**

[1] Research Group for Genomic Epidemiology, Technical University of Denmark, Kgs. Lyngby, Denmark

1   *Correspondence: sjpa@dtu.dk

2   Technical University of Denmark, 2800 Kongens Lyngby, Denmark.

3

## *Abstract*

Plasmids can provide a selective advantage for microorganisms to survive and adapt to new environmental conditions. Plasmid-encoded traits, such as antimicrobial resistance (AMR) or virulence, impact on the ecology and evolution of bacteria and can significantly influence the burden of infectious diseases. Insight about the identity and functions encoded on plasmids on the global scale are largely lacking. Here we investigate the plasmidome of 24 samples (22 countries, 5 continents) from the global sewage surveillance project. We obtained 105 Gbp Oxford Nanopore and 167 Gbp Illumina DNA sequences from plasmid DNA preparations and assembled 165,302 contigs (159,322 circular). Of these, 58,429 encoded for genes with plasmid-related and 11,222 with virus/phage-related proteins. About 90% of the circular DNA elements did not have any similarity to known plasmids. Those that exhibited similarity, had similarity to plasmids whose hosts were previously detected in these sewage samples (e.g. *Acinetobacter, Escherichia, Moraxella, Enterobacter, Bacteroides*, and *Klebsiella*). Some AMR classes were detected at a higher abundance in plasmidomes (e.g. macrolide-lincosamide-streptogramin B, macrolide, and quinolone), as compared to the respective complex sewage samples. In addition to AMR genes, a range of functions were encoded on the candidate plasmids, including plasmid replication and maintenance, mobilization, and conjugation. In summary, we describe a laboratory and bioinformatics workflow for the recovery of plasmids and other potential extrachromosomal DNA elements from complex microbiomes. Moreover, the obtained data could provide further valuable insight into the ecology and evolution of microbiomes, knowledge about AMR transmission, and the discovery of novel functions.

## *Importance*

This is, to the best of our knowledge, the first study to investigate plasmidomes at a global scale using long read sequencing from complex untreated domestic sewage. Previous metagenomic surveys have detected AMR genes in a variety of environments, including sewage. However, it is unknown whether the AMR genes were encoded on the microbial chromosome or are located on extrachromosomal elements, such as plasmids. Using our approach, we recovered a large number of plasmids, of which most appear novel. We identified distinct AMR genes that were preferentially located on plasmids, potentially contributing to their transmissibility. Overall, plasmids are of great importance for the biology of microorganisms in their natural environments (free-living and host-associated), as well as molecular biology, and biotechnology. Plasmidome collections may therefore be valuable resources for the discovery of fundamental biological mechanisms and novel functions useful in a variety of contexts.

2

## *Introduction*

The term plasmid was introduced by Joshua Lederberg in 1952 to describe any extrachromosomal genetic particle (1). It was not until the 1970s when interest in plasmid research rapidly increased and plasmids were introduced as cloning vectors into an area that was dominated by phages as the vector for the transfer of pieces of DNA of choice (2). Since then, plasmids have been highly valuable tools in molecular microbiology. In their natural environment, plasmids are considered key players in horizonal gene transfer. They play crucial roles in the ecology and evolution of bacteria, including their pathogenicity as they can carry virulence factors such as toxins as well as antimicrobial resistances genes (3) (4–6). However, the global diversity of plasmids and distribution of antimicrobial resistance genes are yet to be revealed.

The presence of antimicrobial resistance genes on plasmids are of major interest in the clinical and veterinary areas since they can render prescribed antibiotics for treating pathogens ineffective. There have been a range of large-scale metagenomic-based surveys of antimicrobial resistance genes in soils, humans, animals, plants, and sewage (7–12). However, the genomic context of the AMR genes is largely unknown; for example, whether they are located in the bacterial genome or on plasmids. Such knowledge would be of great value to better assess their potential transmissibility rates and global impact of AMR-gene carrying plasmids on human health.

Plasmids are usually circular DNA elements in bacterial cells, but they can also occur in linear form and be present in archaea and eukaryotic organisms. The size of plasmids is highly variable, ranging from 1,000 bases to hundreds of kilobases. They are present in different quantities (copy numbers) in bacterial cells, varying from a single copy to hundreds of copies in a single cell. This intrinsic and unique nature of plasmids brings about several challenges in plasmidome research (i.e. research on the collective plasmid content in a sample). For example, the low plasmid/chromosome DNA ratio and potential low copy numbers can make it difficult to detect plasmids. These challenges are amplified when plasmidomes are examined from relatively low-cell-density environments such as wastewater. Even assembling and identifying plasmids with low copy number from high biomass samples including single isolates from whole genome sequencing (WGS) data can be challenging. To address these challenges, different approaches have been developed to increase the recovery of plasmids from WGS data (13–16).

Plasmids have now also been recovered from more complex microbiomes using a number of strategies. This includes multiple displacement amplification (MDA) with phi29 DNA polymerase prior to DNA sequencing (17), long read sequencing technology of plasmid DNA, or involvement of advanced assembly strategies (18–21). These studies have however been restricted to a single or few locations, and there is limited knowledge on similarity and differences between plasmids from a

74  range of geographical locations (22–26). We recently showed differences in the AMR gene profiles

75  in urban sewage around the globe, but the location of these AMR genes in the bacteria remains

76  unknown (7).

77  To explore the plasmidome of global sewage, which is characterized by low bacterial cell numbers

78  and challenges to isolate plasmid DNA as previously shown (23–27), we here employ an optimized

79  plasmid DNA isolation procedure, followed by both, plasmid-safe DNase treatment and MDA to

80  obtain sufficient plasmid DNA for Oxford Nanopore sequencing from global urban sewage samples.

81  To improve plasmidome characterizations, we develop an assembly workflow, utilizing the long-read

82  length from the Oxford Nanopore MinION sequencer and Illumina sequences. We obtain thousands

83  of circular candidate plasmid sequences and explore their predicted function.

84

## Material and Methods

### Sample collection and preparation

87  From the global sewage sample collection (7), we selected 24 samples from 22 countries (see Table

88  S1 in the supplementary material). The samples originated from the five most populated continents

89  on Earth and for which we had sufficient sample material available. From each sample, a sewage

90  pellet was collected from 250 ml untreated sewage by centrifugation at 10,000xg for 10 minutes at

91  5°C. The sewage pellets were stored at -80°C until use.

92

### Plasmid DNA extraction and enrichment

94  Plasmid DNA isolation was performed on individual sewage pellets (420 mg) using Plasmid

95  Purification Mini Kit (Qiagen, Cat No./ID: 12123) with QIAGEN-tip 100 (Qiagen, Cat No./ID:

96  10043) following the manufacturer's instruction with the following minor modifications: protein

97  precipitation with P3 buffer mixture was incubated on ice for 15 minutes, elution buffer QF and EB

98  buffer were preheated at 65°C prior to application, and the DNA pellet washing step was done using

99  ice-cold 70% ethanol after isopropanol precipitation. LyseBlue dye for cell lysis indication was

100 added, and all buffer volumes were adjusted to sewage pellet weight. The plasmid DNA pellet was

101 dissolved in 35 µl EB buffer for 1 hour at room temperature. Linear chromosomal DNA was reduced

102 by Plasmid-Safe ATP-Dependent DNase (Epicentre, USA) treatment for 24 hours at 37°C according

103 to the manufacturer's instructions. The DNase was inactivated at 70°C for 30 minutes. To selectively

104 enrich for circular DNA, the plasmid-Safe DNase-treated DNA was amplified using phi29 DNA

105 polymerase (New England Biolabs, USA) following the manufacturer's instructions, similar to as

106 previously described (22). The plasmid DNA is amplified through rolling circle amplification by the

107 phi29 DNA polymerase using random primers, generating multiple DNA replication forks (17). This

4

108   results in long DNA fragments that contain tandem copies (tandem repeats) of the same plasmid.

109   Blank controls were used during plasmid DNA extractions and plasmid enrichment treatments. All

110   negative controls had undetectable DNA measurements using Qubit double-stranded DNA (dsDNA)

111   BR assay kit on a Qubit 2.0 fluorometer (Invitrogen, Carlsbad, CA).

112

### *Plasmid DNA quality assessment*

114   The plasmid DNA yields from the sewage samples were evaluated using gel electrophoresis and

115   Qubit double-stranded DNA (dsDNA) BR assay kit on a Qubit 2.0 fluorometer (Invitrogen, Carlsbad,

116   CA). Plasmid DNA purity was measured and validated by absorbance ratio of 260/280 and 260/230

117   using NanoDrop 100 (ThermoFisher). During pilot experiments that were aimed at protocol

118   development and plasmid DNA enrichment, we also assessed the quality of our plasmid DNA

119   preparations using a 2100 Bioanalyzer (Agilent).

120

### *Library preparation and Oxford Nanopore sequencing*

122   One μg plasmid DNA in 45 μl buffer was used for library preparation. DNA was used without

123   fragmentation. End repair and dA-tailing were performed using NEBNext FFPE Repair Mix (New

124   England BioLabs, 6630) and NEBNext® Ultra™ II End Repair/dA-Tailing Module (New England

125   BioLabs, 7546). DNA was mixed with 3.5 μl NEBNext FFPE DNA Repair Buffer, 2 μl NEBNext

126   FFPE DNA Repair Mix, 3.5 μl Ultra II End-prep reaction buffer and 3 μl Ultra II End-prep enzyme

127   mix and volume was adjusted to 60 μl with nuclease-free water. The reaction tube was flicked 3

128   times and incubated at 20°C for 10 minutes, then inactivated by heating at 65°C for 10 minutes.

129   Clean-up was done using 60 μl Agencourt AMPure XP beads. The beads-reaction suspension was

130   incubated on a HulaMixer at the lowest speed for 10 minutes, followed by two washes with freshly

131   prepared 70% ethanol. DNA was then eluted from the beads in 61 μl 65°C preheated nuclease-free

132   water. A 1 μl DNA aliquot was assessed with Qubit dsDNA BR assay to ensure >700 ng were

133   recovered. A volume of 60 μl of dA-tailed plasmid DNA were added to 25 μl Ligation Buffer (LNB),

134   10 μl NEBNext Quick T4 DNA Ligase NEBNext Quick Ligation Module (New England BioLabs,

135   6056) and 5 μl Adapter Mix (AMX), and mixed by flicking the tube 3-4 times followed by incubation

136   at room temperature for an extended time of 1 hr. The adaptor-ligated plasmid DNA was cleaned up

137   by adding 40 μl Agencourt AMPure XP beads, and the reaction was mixed by flicking the tube and

138   followed by incubation on a HulaMixer at the lowest speed for 10 minutes. The beads were pelleted

139   and resuspended twice in 250 μl Long Fragment Buffer LFB buffer (SQK-LSK109 kit, Oxford

140   Nanopore Technologies). The cleaned adaptor-ligated DNA was then eluted by incubating the pellet

141   in 15 μl Elution Buffer (SQK-LSK109 kit, Oxford Nanopore Technologies) for 20 minutes at room

142   temperature, then transferring the supernatant to a new tube as constructed library. Flowcell priming

143    and library loading preparation were performed according to the manufacturer's instruction (SQK-

144    LSK109 kit, Oxford Nanopore Technologies). Libraries were loaded on FLO-MIN106 R 9.4.1

145    Oxford Nanopore flowcells, and sequencing was run for 48 hours with MinKNOW software default

146    settings.

147

148    *Illumina Sequencing*

149    The enriched plasmid DNA samples were also subjected to Illumina NextSeq sequencing for

150    downstream error-correction of contigs. Libraries were prepared using Nextera XT DNA Library

151    Preparation Kit (Illumina, USA) following the manufacturer's instructions. The libraries were

152    sequenced using NextSeq 550 system (Illumina) with 2 X 150 bp paired-end sequencing per flow

153    cell.

154

155    *Data processing*

156    Basecalling of Nanopore reads was performed using the guppy basecaller (version 3.0.3+7e7b7d0)

157    with the dna_r9.4.1_450bps_hac (high accuracy) configuration. Adapter trimming was performed

158    using porechop (version 0.2.3) downloaded from https://github.com/rrwick/Porechop using the

159    default parameters. Illumina sequencing data were quality and adapter trimmed using bbduk from the

160    bbmap suite (https://sourceforge.net/projects/bbmap/, version 38.23) using the following settings:

161    qin=auto, k=19, rref=adapters.txt, mink=11, qtrim=r, trimq=20, minlength=50, tbo, ziplevel=6,

162    overwrite=t, statscolumns=5.

163

164    *Plasmid assembly from single Nanopore reads*

165    Nanopore reads shorter than 10,000 bases were discarded. Each remaining read was cut into 1,500

166    bases long fragments and passed to the assembly step. The initial fragmentation step of the reads is

167    needed since each read, amplified from a circular element during sample preparation, consists of

168    multiple tandem repeats of the circular element. This is done to eliminate the tandem repeats as well

169    as increase the accuracy of the resulting candidate plasmid DNA sequence. We set the cutting

170    threshold to 1.5 kb to balance between preserving the benefits of long read sequencing and

171    accounting for the error rate of Nanopore sequencing. We decided for a length threshold for cutting

172    (i.e. 1.5 kbp) to not create candidate plasmid DNA sequences from small plasmids that contain

173    multiple copies of the same plasmid. We set the cutting threshold to 1.5 kbp to balance between

174    preserving the benefits of long read sequencing and the error rate of Nanopore sequencing. We also

175    preferred to keeping the cutting threshold more towards the short range to not create candidate

176    plasmids form small plasmids that contain multiple copies of the same plasmid sequence. Read

177 fragments originating from one single read were assembled using minimap2 (version 2.17-r943-dirty)

178 in combination with miniasm version 0.3-r179 (parameter -s 800bp), and error corrected using racon

179 version 1.3.3 (28–30). Assembled contigs were discarded if, after mapping the assembled contig back

180 to the original Nanopore read, hits did not span more than 60% of the read, and if two hits overlapped

181 by more than 50 bp. Assembled candidate contigs were error-corrected using 5 iterations of pilon

182 version 1.23 using the respective Illumina reads from the same sample (31). Candidate contigs longer

183 than 1,000 bases were used for downstream analyses. A schematic overview of the method is

184 presented in Figure 1A.

185

186 *Global plasmidome analysis*

187 To examine the obtained plasmids from our global sewage collection in relation to already known

188 plasmids, we compared our obtained candidate plasmid DNA sequences to the DNA sequences in the

189 plasmid database (PLSDB) using the webtool of PLSDB version 2019_10_07 (32). We used search

190 strategy 'Mash screen' with a maximum p-value of 0.005 and minimum identity of 95%, as well as

191 the option 'winner-takes-all strategy. Samples with less than 100 circular assembled contigs were

192 removed from the analysis as well as genera with less than 10 occurrences over all samples. A

193 clustering of samples was performed using Euclidean distance of the clr-transformed values.

194 Furthermore, all candidate plasmid sequences were sketched using MASH version 2.2 (33). The

195 MASH-distances between all samples were calculated using default settings, resulting in a 24 by 24

196 distance table that was used for principal component analysis.

197

198 *Antimicrobial resistance gene detection analysis*

199 The trimmed Nanopore and Illumina reads were mapped against the ResFinder database (2020-01-

200 25) using kma (version 1.3.0) (34, 35). The Nanopore reads were mapped with settings: mem_mode,

201 ef, nf, bcNano, and bc=0.7. Illumina reads were mapped with settings: mem_mode, ef, nf, 1t1, cge,

202 and t=1. Resistance genes were counted across variants, for example the alleles tet(A)_4_AJ517790

203 and tet(A)_6_AF534183 were both counted as tet(A). Centered log ratios were calculated using the

204 pyCoDa package (https://bitbucket.org/genomicepidemiology/pycoda/src/master/).

205

206 *Gene prediction and functional analysis*

207 Gene prediction was performed using prodigal version 2.6.3, and annotation of protein families was

208 done using hmmscan from HMMER3 version 3.3.1 (http://hmmer.org/) against the pfam database

209 version 33 (36, 37). Predicted genes as well as functional annotation were rejected if the p-value was

210 above 0.000001. Gene ontology (GO) annotations for Pfam IDs were acquired using the mapping of

211 Pfam entries to GO terms as described by Mitchell *et al*. (38).

7

212  To distinguish between potential plasmid and non-plasmid contigs, we used a scheme described

213  previously (39). The scheme contains Pfam identifiers highly specific for plasmids and viruses.

214  Proteins with a plasmid replication initiator protein Rep_3 (PF01051) domain (n=24,824) were

215  investigated further together with the full set of reference Rep_3-domain proteins (n=1,637)

216  downloaded from Pfam (version 33.1). The two data sets were combined and Rep_3-domain proteins

217  with a length of <40 aa residues were discarded, resulting in a data set of 16,930 Rep_3 (PF01051)

218  domain proteins. The protein sequences were aligned using MAFFT (version 7.221) as part of the

219  Galaxy platform (40, 41). A phylogenetic tree was then build using FastTree (version 2.1.10) (42)

220  and visualized using FigTree (version 1.4.4) (https://github.com/rambaut/figtree/releases).

221

## *Generation of plasmid maps*

223  The 50 longest assemblies from each sample were annotated using Prokka (43). Contigs of interest

224  were chosen for mapping based on the presence of known plasmid-encoded genes, such as replication

225  and mobilization systems, toxin-antitoxin pairs, and AMR genes. Plasmids were inspected and

226  visualized using DNAPlotter (44) and Geneious Prime version 2020.2.4 (www.geneious.com). If a

227  coding sequence (CDS) from the Prokka analysis remained unannotated, it was manually annotated

228  by using BLAST search function against the nr database (45) and scanned with HMMER3 against the

229  Pfam database as described above.

230

231

## Results

### Nanopore and Illumina sequencing ouput from plasmid DNA-enriched global sewage samples

The sequencing of 24 plasmid-enriched DNA preparations from untreated sewage from 5 continents (Africa, Asia, Europe, North America, and South America) using Oxford Nanopore sequencing technology produced 1.2 to 9.7 Gbp (median 3.5 Gbp) sequencing data per sample (see Table S1 in the supplementary material). The median read length was 7.3 kb (range 1,075 to 11,018 bases) (see Figure S1 in the supplemental material). After quality trimming and removing sequences below 10,000 bases, the median sequencing throughput was 1.9 Gbp and the median read length 23,000 bases (see Table A at https://doi.org/10.6084/m9.figshare.13395446). The Illumina generated sequencing data per sample were between 1.5 and 9.7 Gbp with a median of 4.8 Gbp after adapter and quality trimming. A median of 41 million paired-end reads per sample were obtained (see Table B at https://doi.org/10.6084/m9.figshare.13395446).

### Circular DNA sequences obtained using single Oxford Nanopore reads

Upon assembly and polishing (Figure 1A), we obtained a total of 165,302 contigs from the 24 sewage samples, of which 159,322 contigs (96.4%) were suggested by miniasm to be circular (Figure 1B, and see Table C at https://doi.org/10.6084/m9.figshare.13395446). The longest assembled circular contig had a length of 17.4 kbp and was obtained from a sample in Brasil (BRA.1, South America). Most of the circular contigs were obtained from the Tanzanian (TZA, Africa) sewage sample, and they had an average length of 1.7 kbp (see Table C at https://doi.org/10.6084/m9.figshare.13395446).

### Classification of assembled circular DNA elements

To obtain information about the identity of the obtained circular DNA elements, we performed gene prediction, annotation, and classification based on plasmid- and virus/phage-specific Pfam domains (39). Overall, we detected Pfam domains (including domains of unknown function (DUF)) on 47.01% of the circular elements, potentially suggesting the presence of many novel DNA sequences not encoding for known protein domains. For the DNA elements (circular & linear) for which Pfam domains were detected, the majority (88.39%) contained predicted genes with plasmid- or virus/phage-related Pfam entries (see Figure 2, Figure S2 in the supplementary material, and Table D at https://doi.org/10.6084/m9.figshare.13395446). Overall, we found 55,337 circular DNA elements that encoded for known plasmid-related Pfam domains (and not viral-related Pfam domains). The highest number of plasmid-related candidate sequences were detected in the sample from the Czech

266    Republic (CZE, Europe), followed by Tanzania (TZA, Africa), and Kosovo (XK, Europe). The

267    sample from China (CHN, Asia) was the only sample from which more potential virus/phage-related

268    contigs than candidate plasmids were obtained (see Figure 2, Figure S2 in the supplementary

269    material, and Table D at https://doi.org/10.6084/m9.figshare.13395446).

270    On the circular elements with plasmid-related Pfam domains, protein families involved in plasmid

271    replication were the most abundant and they included Relaxase, *Rep_1, Rep_2, Rep_3, Rep_trans,*

272    *RepL, and Replicase* (Figure 2A). For example, we detected a total of 24,824 open reading frames

273    with a plasmid replication initiator protein Rep_3 (PF01051) domain. Even though Rep_3-domain

274    proteins from all continents were observed across the phylogenetic tree, some clades mainly

275    represented proteins from one continent, interspersed with protein sequences from other continents

276    (Figure 2B). For instance, clades that mainly harbored proteins originating from Europe, also

277    frequently contained protein sequences from North America and other continents. Clades dominated

278    by Rep_3 (PF01051) domain proteins from Africa also frequently harbored similar proteins from

279    South America.

280    Furthermore, protein families involved in plasmid mobilization were detected, such as Mob_Pre,

281    *MobA_MobL,* and *MobC* (Figure 2A). In addition, we identified protein families related to

282    virus/phage replication and capsid proteins, as well as protein domains binding to DNA (HTH_17,

283    HTH_23, HTH_Crp_2) and that might be involved in regulating gene expression.

284

285    *Global plasmidome pattern based on known plasmids*

286    To examine whether our collection of plasmid sequences contained already known sequences, we

287    compared the obtained plasmid DNA sequences to the entries in the plasmid database (PLSDB). This

288    analysis revealed that only 10.1% of our circular elements were similar to known plasmids (see Table

289    E at https://doi.org/10.6084/m9.figshare.13395446). The majority of plasmids that exhibited some

290    similarity to entries in the PLDB originated from *Acinetobacter* (33%), *Enterococcus* (21%) as well

291    as *Flavobacterium* (10%); genera that were previously detected in these sewage microbiomes (7).

292    Overall, most plasmids with similarities to already known ones were found in the samples from India,

293    Kosovo, Pakistan, Czech Republic, Iceland, and Brazil (see Table E at

294    https://doi.org/10.6084/m9.figshare.13395446). Clustering analysis of the abundancies of plasmids

295    with known relatives in PLSDB revealed three main clusters (Figure 3A). The first cluster comprised

296    samples that overall exhibited a low number of known plasmids and included samples from Europe

297    (ALB, POL, ESP, SVN) and a sample from Ghana. The second cluster included samples with

298    plasmids from a large range of bacterial genera at higher abundance, and comprised samples from

299    Europe (ISL, DEU, CZE), North America (USA.1, USA.2, CAN), India, Brazil and Tanzania. The

300    third cluster comprised samples with known plasmids from few bacterial genera and included

301  samples from Asia (CHN, PAK), Africa (CIV), Europe (XK), and South America (ECU, PER)

302  (Figure 3A).

303  In a principal component analysis of the same data, a similar clustering was observed. Furthermore,

304  along the first principal component, samples from Asia and Europe appeared to be most different

305  from each other and with samples from Africa, and North and South America in between. Upon

306  examining the particular reference plasmids and their bacterial hosts that were driving this pattern a

307  similar observation was made: plasmids from bacterial hosts originating from Europe appeared to

308  segregate along the first principle component from plasmids and their bacterial hosts originating from

309  Asia (Figure 3B). This observation was supported by a cluster analysis on plasmid-level, in which

310  five clusters were observed: Samples from Europe did not cluster with samples from Asia, and

311  different sets of known plasmids were found in the samples from Europe and Asia, respectively (see

312  Figure S3 in the supplemental material). Generally, only few known plasmids were detected in the

313  samples from Albania, Slovenia, Spain, Poland, Ecuador, and Ghana (see Figure S3 in the

314  supplementary material and Table E at https://doi.org/10.6084/m9.figshare.13395446).

315  Given the large fraction of candidate plasmid sequences that did not exhibit similarity to already

316  known plasmids, we performed a reference-independent analysis by calculating Mash-distances

317  based on all plasmid sequences for each sample. In this analysis, the plasmidomes clustered in two

318  main clusters (see Figure S4 in the supplemental material). The first cluster harbored all samples

319  from Europe (with the exception of Poland), as well as the samples from Canada (North America),

320  Pakistan and India (Asia), and Côte d'Ivoire (Africa). The second cluster harbored all samples from

321  South America, both samples from the USA (North America), as well as Tanzania and Ghana

322  (Africa), and China (Asia) (see Figure S4 in the supplemental material). This suggests that the

323  sequence space encompassing novel plasmid sequences (i.e. those that did not exhibit similarity to

324  sequences in the PLSDB) provides an extended, yet to be discovered, dimension into plasmid

325  ecology and evolution.

326

327  *Antimicrobial resistance genes in plasmidomes*

328  To gain insight into antimicrobial resistance genes on the plasmids from sewage, and compare them

329  to those detected in the whole community of the same sewage samples, we performed a ResFinder

330  analysis on three sequencing read data sets: whole community DNA sequenced using Illumina (7),

331  plasmidome DNA sequenced using Illumina (this study), and plasmidome DNA sequenced using

332  Nanopore sequencing (this study).

333  Overall, many of the antimicrobial resistance genes and antimicrobial classes that were detected

334  using whole community sequencing, were also detected in the two plasmidome datasets, with a few

335  exceptions. For example, the two antimicrobial classes macrolide-streptogramin B and lincosamide-

336    pleuromutilin-streptogramin A were not detected in the plasmidome samples in about half of the

337    cases (Figure 4A, and see Figure S5A in the supplementary material, and Tables F and G at

338    https://doi.org/10.6084/m9.figshare.13395446). Occasionally, also genes conferring resistance to

339    other antimicrobial classes were not detected in individual plasmidome samples as compared to the

340    whole community, and these included genes conferring resistance to lincosamide, phenicol, or

341    aminoglycoside. It may be that genes that were detected more frequently in the whole community

342    sample, as compared to the plasmidome samples, are preferentially encoded on the bacterial

343    chromosomes or larger plasmids.

344    Conversely, genes conferring resistance to the antimicrobial classes macrolide-lincosamide-

345    streptogramin B, as well as macrolide, and quinolone, were more frequently observed in the

346    plasmidome samples (Figure 4A, and see Figure S5 in the supplementary material and Tables F and

347    G at https://doi.org/10.6084/m9.figshare.13395446). The most frequently observed AMR genes

348    related to these three classes were *ermB*, *ermT*, *ermF* (macrolide-lincosamide-streptogramin B),

349    *mphE*, *mefA*, *msrD* (Macrolide), and *qnrB19*, *qnrD1*, *qnrD2*, *qnrD3*, *qnrVC4* (Quinolone). The

350    higher frequency of those genes in the plasmidome samples may suggest that they are more

351    frequently found on plasmids in general, or on smaller plasmids as compared to large ones. Another

352    gene that was frequently observed across samples is *msrE*, and which was slightly higher abundant in

353    plasmidomes (average abundance 15.4%, SEM 1.86) as compared to whole community samples

354    (average abundance 11.5%, SEM 1.88). As examples, a few randomly chosen candidate plasmids and

355    their encoded genes, including AMR genes, are displayed in Figure 4B.

356

357    *Functional characterization of plasmidomes*

358    To gain further insight into the functions encoded on all circular elements, we obtained GO

359    annotations for the predicted proteins through mapping of pfam entries to GO terms. A clustering

360    analysis revealed the separation of plasmidomes into two main clusters (see Figure S6 in the

361    supplemental material). Cluster 1 comprised samples from Europe (ISL, CZE, XK, DEU) as well as

362    North America (USA.1, CAN) and South America (BRA.1, ECU). Cluster 2 comprised the samples

363    from Asia (IND, PAK, CHN), Africa (TZA, CIV) and the remaining samples from Europe (POL,

364    ESP, SVN) and South America (PER). This clustering based on protein functions appeared to have

365    some similarity to the clustering based on nucleotide sequence similarity to known plasmids (Figure

366    3). In both analyses, the European samples from ISL, CZE, and DEU exhibited similarities, while the

367    other European samples from POL, ESP, SVN clustered together separately. Furthermore, in both

368    analyses, samples from North America (USA.1, CAN) and South America (BRA.1) clustered with

369    the European samples from ISL, CZE, and DEU.

370 Functions that appeared to be enriched in samples from cluster 1 include, conjugation, recombinase

371 activity, DNA methylation, protein secretion (type IV secretion system), response to antibiotic, toxic

372 substance binding, response to toxic substance, unidirectional conjugation, and bacteriocin immunity

373 (see Figure S5 in the supplemental material). Cluster 2 appeared to overall have fewer proteins that

374 could be annotated using this strategy, and the samples exhibited a higher diversity of functional

375 patterns compared to samples from cluster 1. Some samples from cluster 2 exhibited an enrichment

376 of proteins that may be related to viruses/phages, such as viral capsids, structural molecule activity,

377 RNA binding, RNA helicase activity, and these were in particular samples that appeared to have a

378 higher abundance of virus/phage related Pfam domains (Figure 2). The majority of samples in both

379 clusters harbored proteins involved in plasmid maintenance (see Figure S6 in the supplemental

380 material).

381

## 382 *Discussion*

383 This is the first study to investigate plasmidomes at a global scale using long read sequencing from

384 sewage. We show that our approach facilitated the recovery of complete plasmids from complex

385 metagenomic samples with a sufficient quality to perform gene prediction and functional annotation.

386 In total, we obtained 165,302 DNA elements of which 159,322 were circular. The average length was

387 1.9 kb (min 1 kbp, max 17.4 kbp), suggesting that mainly small plasmids were obtained. This might

388 reflect the true distribution but could also be biased due to a number of reasons, for example, smaller

389 plasmids are more stable and thus have higher chance of getting though the DNA extraction step

390 undamaged. Since a DNase step was used to reduce the amount of chromosomal DNA, damaged

391 plasmids might have been digested as well. Another possibility could be that some plasmids were

392 already damaged during storage and transportation, as the sewage was frozen and shipped, and many

393 of the samples arrived thawed and were frozen again. Another reason could be that our assembly

394 workflow was not able to perform a successful assembly on larger plasmids with a high number of

395 tandem-repeats.

396

397 We identified a range of functions encoded on the candidate plasmids, including plasmid replication

398 and maintenance, mobilization, conjugation, antimicrobial resistance, and bacteriocin immunity.

399 However, not all plasmid-related DNA elements encoded for a plasmid-replication gene, suggesting

400 that they may not be self-replicating DNA molecules. It should though be noted that also already

401 described plasmids do not necessarily encode for a rep gene using current annotation algorithms.

402 Furthermore, we found that about half of the circular DNA elements did not encode for any known

403 Pfam domains. This could suggest that we detected many novel DNA sequences not encoding for

404 known protein domains. A hypothesis could be that a fraction of the circular DNA elements are novel

405 extrachromosomal elements that are hitherto undescribed and may also originate from various

406 domains of life, including bacteria, archaea, and eukaryotes (46–48). Alternatively, open reading

407 frames might not always have been properly detected because of sequencing errors not corrected in

408 the polishing steps with Nanopore and Illumina reads. This could certainly have contributed to it, as

409 we occasionally observed fragmented genes due to remaining sequencing errors, even after polishing.

410 This challenge may be alleviated with the ongoing improvement of Oxford Nanopore chemistry and

411 basecalling algorithms. Nevertheless, collectively, we obtained 58,429 DNA elements (circular &

412 linear) that encoded for proteins with plasmid-related Pfams, and 17,292 circular DNA elements

413 exhibited sequence similarity to known plasmids, suggesting that we successfully discovered many

414 novel candidate plasmid DNA sequences.

415

416 For candidate plasmids that exhibited some similarities to known plasmids, we found that they

417 originated from bacterial taxa previously detected in these complex sewage samples, such as

418 *Acinetobacter, Escherichia, Moraxella, Enterobacter, Bacteroides*, and *Klebsiella* (7). These genera

419 include bacteria that are part of the human gut microbiome and/or opportunistic pathogens. Hence,

420 some of these plasmids might play a role in gut microbial ecology and potential antimicrobial

421 resistance transmission (49, 50). It should be noted, however, that overall, only ~10.1% of our

422 circular elements were similar to known plasmids in the PLSDB, and which may be partly explained

423 by differences in plasmid contents (plasmid average size 1.9 kbp (this study) and 53.2 kbp (PLSDB))

424 (32). In line with this, we observed that the plasmidome samples clustered somewhat differently

425 when all candidate plasmid sequences were taken into account (and not only those that exhibited

426 similarity to known reference plasmids). It will be interesting to investigate our candidate plasmids

427 further in future studies, ideally through involvement of more plasmidome samples and extended

428 metadata. There may be a range of factors that may play role in explaining differences and

429 similarities between plasmidomes, such as climate, population-related differences including human

430 ethnicity, health status, sanitation, and economy including trading between countries.

431

432 Overall, AMR classes that were detected in the plasmidome sequencing data sets were also found in

433 the sequencing data from the whole complex sewage samples, suggesting that the plasmidomes are a

434 good representation of what is present in the complex samples. Some AMR gene classes, however,

435 were more predominant in the whole community (e.g. macrolide-streptogramin B, lincosamide-

436 pleuromutilin-streptogramin A), and others more in the plasmidomes (e.g. macrolide-lincosamide-

437 streptogramin B, macrolide, and quinolone). This could suggest that the AMR genes conferring

438 resistance to the latter AMR gene classes are preferentially located on plasmids as compared to

439    chromosomes. However, given that we mainly recovered small plasmids, it could also be an

440    indication for that the AMR genes preferentially detected in the whole community may be located on

441    large plasmids that were not recovered here. Whether certain abundant AMR genes in the

442    plasmidomes are plasmid- or chromosome-associated may also be dependent on the particular

443    bacterial host (see Figure S7 in the supplemental material) (51).

444

445    While our approach and findings are a significant advancement to previous work, there are still

446    aspects that can be improved in the future. For example, the assembly workflow could be improved

447    to resolve remaining repetitive regions within the plasmid, as a range of circular elements still

448    consisted of tandem-repeats of the actual plasmid sequence. This could potently be solved by

449    introducing a dynamic cutting step using the k-mer composition of the full read. Despite the high

450    error rate of the Nanopore sequencing reads, the raw read should still contain a set of k-mers with 10-

451    15 bases length that could help interfering the appropriate fragmentation length. In addition, the

452    plasmid DNA isolation could be improved significantly to increase a) the overall amount of plasmid

453    DNA (in order to avoid having to perform MDA), and b) the amount of larger plasmids. Further

454    possibilities to identify new plasmids could also involve *in vivo* proximity-ligation Hi-C or single-cell

455    sequencing that would also allow the discovery of new plasmids directly together with their host cell

456    (52, 53).

457

458    Overall, our study provides new insight about the technical applicability of long-read Nanopore

459    sequencing for plasmidome analysis of complex biological samples, as well as a foundation for

460    exploring plasmid ecology and evolution at a global scale. For example, we can now better explore

461    the genomic context of AMR genes, and reveal whether they are located on the microbial

462    chromosome or on mobile genetic elements such as plasmids. This knowledge is of great value in

463    assessing the potential transmissibility of AMR genes with resulting impact on antibiotic treatments

464    in the medical and veterinary sectors and the one health perspective. Furthermore, the dataset

465    provides a valuable resource for further exploring extrachromosomal DNA elements including

466    potential novel functions.

467

468    *Acknowledgment*

15

472    Programme. The funders had no role in study design, data collection and interpretation, or the

473    decision to submit the work for publication.

474    We thank Christina Aaby Svendsen (Technical University of Denmark) for technical support with the

475    Illumina sequencing of the plasmidomes.

476    Sequencing data analysis was performed using the DeiC National Life Science Supercomputer at

477    DTU.

478

479    *Data availability*

480    The DNA sequences generated in this project are available through ENA/GenBank/DDBJ under the

481    accession number PRJEB41171 (Nanopore reads: ERX4715074-ERX4715097; Illumina reads:

482    ERX5299122-ERX5299145; Assemblies: ERZ1694234-ERZ1694257). The code for the creation of

483    assemblies is accessible from Github (https://github.com/philDTU/plasmidPublication) and

484    additional supplementary material is available at

485    https://figshare.com/projects/A_Peek_into_the_Plasmidome_of_Global_Sewage/94448.

486

## *References*

1. Lederberg J. 1952. Cell Genetics and Hereditary Symbiosis. Physiological Reviews 32:403–430.

2. Cohen SN, Chang ACY, Boyer HW, Helling RB. 1973. Construction of Biologically Functional Bacterial Plasmids In Vitro. PNAS 70:3240–3244.

3. Rodríguez-Beltrán J, DelaFuente J, León-Sampedro R, MacLean RC, San Millán Á. 2021. Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. Nature Reviews Microbiology https://doi.org/10.1038/s41579-020-00497-1.

4. Johnson TJ, Nolan LK. 2009. Pathogenomics of the Virulence Plasmids of Escherichia coli. Microbiol Mol Biol Rev 73:750–774.

5. Bratu S, Brooks S, Burney S, Kochar S, Gupta J, Landman D, Quale J. 2007. Detection and Spread of Escherichia coli Possessing the Plasmid-Borne Carbapenemase KPC-2 in Brooklyn, New York. Clin Infect Dis 44:972–975.

6. Tian G-B, Doi Y, Shen J, Walsh TR, Wang Y, Zhang R, Huang X. 2017. MCR-1-producing Klebsiella pneumoniae outbreak in China. The Lancet Infectious Diseases 17:577.

7. Hendriksen RS, Munk P, Njage P, Bunnik B, McNally L, Lukjancenko O, Röder T, Nieuwenhuijse D, Pedersen SK, Kjeldgaard J, Kaas RS, Clausen PTLC, Vogt JK, Leekitcharoenphon P, Schans MGM, Zuidema T, Husman AMR, Rasmussen S, Petersen B, Bego A, Rees C, Cassar S, Coventry K, Collignon P, Allerberger F, Rahube TO, Oliveira G, Ivanov I, Vuthy Y, Sopheak T, Yost CK, Ke C, Zheng H, Baisheng L, Jiao X, Donado-Godoy P, Coulibaly KJ, Jergović M, Hrenovic J, Karpíšková R, Villacis JE, Legesse M, Eguale T, Heikinheimo A, Malania L, Nitsche A, Brinkmann A, Saba CKS, Kocsis B, Solymosi N, Thorsteinsdottir TR, Hatha AM, Alebouyeh M, Morris D, Cormican M, O'Connor L, Moran-Gilad J, Alba P, Battisti A, Shakenova Z, Kiiyukia C, Ng'eno E, Raka L, Avsejenko J, Bērziņš A, Bartkevics V, Penny C, Rajandas H, Parimannan S, Haber MV, Pal P, Jeunen G-J, Gemmell N, Fashae K, Holmstad R, Hasan R, Shakoor S, Rojas MLZ, Wasyl D, Bosevska G, Kochubovski M, Radu C, Gassama A, Radosavljevic V,

512  Wuertz S, Zuniga-Montanez R, Tay MYF, Gavačová D, Pastuchova K, Truska P, Trkov M, Esterhuyse K,

513  Keddy K, Cerdà-Cuéllar M, Pathirage S, Norrgren L, Örn S, Larsson DGJ, Van der Heijden T, Kumburu HH,

514  Sanneh B, Bidjada P, Njanpop-Lafourcade B-M, Nikiema-Pessinaba SC, Levent B, Meschke JS, Beck NK,

515  Van CD, Do Phuc N, Tran DMN, Kwenda G, Tabo D, Wester AL, Cuadros-Orellana S, Amid C, Cochrane G,

516  Sicheritz-Ponten T, Schmitt H, Alvarez JRM, Aidara-Kane A, Pamp SJ, Lund O, Hald T, Woolhouse M,

517  Koopmans MP, Vigre H, Petersen TN, Aarestrup FM. 2019. Global monitoring of antimicrobial resistance

518  based on metagenomics analyses of urban sewage. Nature Communications 10:1124.

519 8. Munk P, Knudsen BE x000E6 r, Lukjacenko O, Duarte ASR, Gompel L, Luiken REC, Smit LAM, Schmitt H,

520  Garcia AD, Hansen RB, Petersen TN, Bossers A, x000E9 ER, Graveland H, van Essen A, Gonzalez-Zorn B,

521  Moyano G, Sanders P, Chauvin C, David J, Battisti A, Caprioli A, Dewulf J, Blaha T, Wadepohl K, Brandt

522  M, Wasyl D, ska MS x00144, Zajac M, Daskalov H, Saatkamp HW, rk KDCS x000E4, Lund O, Hald T, Pamp

523  S x000FC nje J, Vigre H x000E5 kan, Heederik D, Wagenaar JA, Mevius D, Aarestrup FM. 2018.

524  Abundance and diversity of the faecal resistome in slaughter pigs and broilers in nine European

525  countries. Nature Microbiology 1–14.

526 9. Campbell TP, Sun X, Patel VH, Sanz C, Morgan D, Dantas G. 2020. The microbiome and resistome of

527  chimpanzees, gorillas, and humans across host lifestyle and geography. 6. The ISME Journal 14:1584–

528  1599.

529 10. Chen Q-L, Cui H-L, Su J-Q, Penuelas J, Zhu Y-G. 2019. Antibiotic Resistomes in Plant Microbiomes.

530  Trends in Plant Science 24:530–541.

531 11. Forsberg KJ, Patel S, Gibson MK, Lauber CL, Knight R, Fierer N, Dantas G. 2014. Bacterial phylogeny

532  structures soil resistomes across habitats. 7502. Nature 509:612–616.

533 12. Carr VR, Witherden EA, Lee S, Shoaie S, Mullany P, Proctor GB, Gomez-Cabrero D, Moyes DL. 2020.

534  Abundance and diversity of resistomes differ between healthy human oral cavities and gut. 1. Nature

535  Communications 11:693.

536    13.    Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Unicycler: Resolving bacterial genome assemblies from

537            short and long sequencing reads. PLOS Computational Biology 13:e1005595.

538    14.    Antipov D, Hartwick N, Shen M, Raiko M, Lapidus A, Pevzner PA. 2016. plasmidSPAdes: assembling

539            plasmids from whole genome sequencing data. Bioinformatics 32:3380–3387.

540    15.    Vielva L, de Toro M, Lanza VF, de la Cruz F. 2017. PLACNETw: a web-based tool for plasmid

541            reconstruction from bacterial genomes. Bioinformatics 33:3796–3798.

542    16.    Rozov R, Brown Kav A, Bogumil D, Shterzer N, Halperin E, Mizrahi I, Shamir R. 2017. Recycler: an

543            algorithm for detecting plasmids from de novo assembly graphs. Bioinformatics 33:475–482.

544    17.    Dean FB, Nelson JR, Giesler TL, Lasken RS. 2001. Rapid Amplification of Plasmid and Phage DNA Using

545            Phi29 DNA Polymerase and Multiply-Primed Rolling Circle Amplification. Genome Res 11:1095–1099.

546    18.    Che Y, Xia Y, Liu L, Li A-D, Yang Y, Zhang T. 2019. Mobile antibiotic resistome in wastewater treatment

547            plants revealed by Nanopore metagenomic sequencing 1–13.

548    19.    Bertrand D, Shaw J, Kalathiyappan M, Ng AHQ, Kumar MS, Li C, Dvornicic M, Soldo JP, Koh JY, Tong C,

549            Ng OT, Barkham T, Young B, Marimuthu K, Chng KR, Sikic M, Nagarajan N. 2019. Hybrid metagenomic

550            assembly enables high-resolution analysis of resistance determinants and mobile elements in human

551            microbiomes. Nature Biotechnology 1–15.

552    20.    Antipov D, Raiko M, Lapidus A, Pevzner PA. 2019. Plasmid detection and assembly in genomic and

553            metagenomic datasets. Genome Res gr.241299.118.

554    21.    Jørgensen TS, Hansen MA, Xu Z, Tabak MA, Sørensen SJ, Hansen LH. 2017. Plasmids, Viruses, And Other

555            Circular Elements In Rat Gut. bioRxiv 143420.

556    22.    Kav AB, Sasson G, Jami E, Doron-Faigenboim A, Benhar I, Mizrahi I. 2012. Insights into the bovine rumen

557            plasmidome. PNAS 109:5452–5457.

558   23.  Kav AB, Rozov R, Bogumil D, Sørensen SJ, Hansen LH, Benhar I, Halperin E, Shamir R, Mizrahi I. 2020. Unravelling plasmidome distribution and interaction with its hosting microbiome. Environmental Microbiology 22:32–44.

561   24.  Zhang T, Zhang X-X, Ye L. 2011. Plasmid Metagenome Reveals High Levels of Antibiotic Resistance Genes and Mobile Genetic Elements in Activated Sludge. PLOS ONE 6:e26041.

563   25.  Sentchilo V, Mayer AP, Guy L, Miyazaki R, Green Tringe S, Barry K, Malfatti S, Goessmann A, Robinson-Rechavi M, van der Meer JR. 2013. Community-wide plasmid gene mobilization and selection. 6. The ISME Journal 7:1173–1186.

566   26.  Kothari A, Wu Y-W, Chandonia J-M, Charrier M, Rajeev L, Rocha AM, Joyner DC, Hazen TC, Singer SW, Mukhopadhyay A. 2019. Large Circular Plasmids from Groundwater Plasmidomes Span Multiple Incompatibility Groups and Are Enriched in Multimetal Resistance Genes. mBio 10.

569   27.  Kav AB, Sasson G, Jami E, Doron-Faigenboim A, Benhar I, Mizrahi I. 2012. Insights into the bovine rumen plasmidome. PNAS 109:5452–5457.

571   28.  Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34:3094–3100.

572   29.  Li H. 2016. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. Bioinformatics 32:2103–2110.

574   30.  Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo genome assembly from long uncorrected reads. Genome Res 27:737–746.

576   31.  Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J, Young SK, Earl AM. 2014. Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. PLOS ONE 9:e112963.

579   32.  Galata V, Fehlmann T, Backes C, Keller A. 2019. PLSDB: a resource of complete bacterial plasmids. Nucleic Acids Res 47:D195–D202.

581  33.  Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM. 2016. Mash: fast

582  genome and metagenome distance estimation using MinHash. Genome Biology 17:132.

583  34.  Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, Aarestrup FM, Larsen MV.

584  2012. Identification of acquired antimicrobial resistance genes. Journal of Antimicrobial Chemotherapy

585  67:2640–2644.

586  35.  Clausen PTLC, Aarestrup FM, Lund O. 2018. Rapid and precise alignment of raw reads against

587  redundant databases with KMA. BMC Bioinformatics 19:307.

588  36.  Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal: prokaryotic gene

589  recognition and translation initiation site identification. BMC Bioinformatics 11:11:119.

590  37.  El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, Qureshi M, Richardson LJ, Salazar GA,

591  Smart A, Sonnhammer ELL, Hirsh L, Paladin L, Piovesan D, Tosatto SCE, Finn RD. 2019. The Pfam protein

592  families database in 2019. Nucleic Acids Res 47:D427–D432.

593  38.  Mitchell A, Chang H-Y, Daugherty L, Fraser M, Hunter S, Lopez R, McAnulla C, McMenamin C, Nuka G,

594  Pesseat S, Sangrador-Vegas A, Scheremetjew M, Rato C, Yong S-Y, Bateman A, Punta M, Attwood TK,

595  Sigrist CJA, Redaschi N, Rivoire C, Xenarios I, Kahn D, Guyot D, Bork P, Letunic I, Gough J, Oates M, Haft

596  D, Huang H, Natale DA, Wu CH, Orengo C, Sillitoe I, Mi H, Thomas PD, Finn RD. 2015. The InterPro

597  protein families database: the classification resource after 15 years. Nucleic Acids Res 43:D213–D221.

598  39.  Jørgensen TS, Hansen MA, Xu Z, Tabak MA, Sørensen SJ, Hansen LH. 2017. Plasmids, Viruses, And Other

599  Circular Elements In Rat Gut. bioRxiv 143420.

600  40.  Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7: Improvements

601  in Performance and Usability. Molecular Biology and Evolution 30:772–780.

602  41.  Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Čech M, Chilton J, Clements D, Coraor N, Grüning

603  BA, Guerler A, Hillman-Jackson J, Hiltemann S, Jalili V, Rasche H, Soranzo N, Goecks J, Taylor J,

604      Nekrutenko A, Blankenberg D. 2018. The Galaxy platform for accessible, reproducible and collaborative
605      biomedical analyses: 2018 update. Nucleic Acids Research 46:W537–W544.

606  42.  Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees for Large
607      Alignments. PLOS ONE 5:e9490.

608  43.  Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics 30:2068–2069.

609  44.  Carver T, Thomson N, Bleasby A, Berriman M, Parkhill J. 2009. DNAPlotter: circular and linear
610      interactive genome visualization. Bioinformatics 25:119–120.

611  45.  Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. Journal of
612      molecular biology 215:403–410.

613  46.  Lanciano S, Carpentier M-C, Llauro C, Jobet E, Robakowska-Hyzorek D, Lasserre E, Ghesquière A,
614      Panaud O, Mirouze M. 2017. Sequencing the extrachromosomal circular mobilome reveals
615      retrotransposon activity in plants. PLOS Genetics 13:e1006630.

616  47.  Shibata Y, Kumar P, Layer R, Willcox S, Gagan JR, Griffith JD, Dutta A. 2012. Extrachromosomal
617      microDNAs and chromosomal microdeletions in normal tissues. Science 336:82–86.

618  48.  Møller HD, Mohiyuddin M, Prada-Luengo I, Sailani MR, Halling JF, Plomgaard P, Maretty L, Hansen AJ,
619      Snyder MP, Pilegaard H, Lam HYK, Regenberg B. 2018. Circular DNA elements of chromosomal origin
620      are common in healthy human somatic tissue. 1. Nature Communications 9:1069.

621  49.  San Millan A. 2018. Evolution of Plasmid-Mediated Antibiotic Resistance in the Clinical Context. Trends
622      in Microbiology 26:978–985.

623  50.  Ogilvie LA, Firouzmand S, Jones BV. 2012. Evolutionary, ecological and biotechnological perspectives on
624      plasmids resident in the human gut mobile metagenome. Bioengineered 3:13–31.

625    51.  Alcock BP, Raphenya AR, Lau TTY, Tsang KK, Bouchard M, Edalatmand A, Huynh W, Nguyen A-LV, Cheng

626          AA, Liu S, Min SY, Miroshnichenko A, Tran H-K, Werfalli RE, Nasir JA, Oloni M, Speicher DJ, Florescu A,

627          Singh B, Faltyn M, Hernandez-Koutoucheva A, Sharma AN, Bordeleau E, Pawlowski AC, Zubyk HL,

628          Dooley D, Griffiths E, Maguire F, Winsor GL, Beiko RG, Brinkman FSL, Hsiao WWL, Domselaar GV,

629          McArthur AG. 2020. CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic

630          resistance database. Nucleic Acids Res 48:D517–D525.

631    52.  Stalder T, Press MO, Sullivan S, Liachko I, Top EM. 2019. Linking the resistome and plasmidome to the

632          microbiome. 10. The ISME Journal 13:2437–2446.

633    53.  Lan F, Demaree B, Ahmed N, Abate AR. 2017. Single-cell genome sequencing at ultra-high- throughput

634          with microfluidic droplet barcoding. Nature Biotechnology 35:640–646.

635

636

637     *__Figure legends__*

638

639     **Figure 1. Schematic overview of the single read assembly workflow and size distribution of**

640     **assembled reads.** A) Nanopore reads (based on plasmid DNA amplified with phi29) longer than

641     10,000 bases were split into 1,500 bases long fragments. The sequence fragments were then

642     assembled using minimap2 and miniasm and subsequently polished two times: 1. with the Nanopore

643     fragments using racon and 2. with the Illumina reads using pilon. B) The size distribution of circular

644     (orange) and linear (violet) assembled elements. These are the candidate plasmid sequences that

645     successfully mapped to the original Nanopore read (i.e. covering more than 60% of the read, and not

646     overlapping by more than 50 bp for multiple hits). Of the total 165,302 assemblies, 159,322 were

647     characterized to be circular and 5,980 to be linear.

648

649     **Figure 2. Functional characterization of circular DNA elements based on protein families**. A)

650     The bar plot displays the fraction of Pfam identifiers assigned to predicted proteins on the circular

651     elements. The 31 Pfam identifiers represent the Top10 Pfam identifiers for each sample. Protein

652     domains specifically involved in plasmid mobilization and plasmid replication are indicated by red

653     and blue colors, respectively (see legend to the bottom right). Virus/phage related Pfam identifiers

654     are indicated in green colors. Remaining Pfam identifiers are grouped (other) and indicated by dark

655     grey. B) The dataset of proteins with a Rep_3 (PF01051) domain (n= 24,824) were combined

656     together with the 1,637 reference Rep_3 (PF01051) proteins from Pfam. The protein sequences with

657     a length of $>/= 40$ aa (n=16,930) were aligned using MAFFT. A phylogenetic tree was build using

658     FastTree and visualized using FigTree. A high-resolution version of the phylogenetic tree is available

659     from Figshare at https://doi.org/10.6084/m9.figshare.14112992.

660

661     **Figure 3. Comparison of candidate plasmids from global sewage with known plasmids in**

662     **plasmid database (PLSDB).** A) Heat map of centered log ratio (clr)-transformed abundancies of

663     plasmid candidates assigned to plasmids in the PLSDB at bacterial genus level. The phylum level is

664     indicated in parenthesis, A: Actinobacteria; B: Bacteroidetes; aP: alpha-Proteobacteria; bP: beta-

665     Proteobacteria; gP: gamma-Proteobacteria; F: Firmicutes. Clustering of samples was performed using

666     Euclidean distance of the clr-transformed values. B) Principal component analysis of clr-transformed

667     abundancies of known plasmids detected by the PLSDB. The plot on the top reveals similarities and

668     differences between samples. The plot in the bottom reveals the known plasmids that drive the

669     partitioning of the samples, with 17.6% of the variation explained by the first and 11.1% by the

670     second principal component.

671

672 **Figure 4. Antimicrobial resistance profiles from the whole community and plasmidomes from**

673 **global sewage. A)** Bar plot displaying the proportions of antimicrobial resistance classes detected in

674 a ResFinder-based analysis using the Illumina reads from the whole community, as well as Illumina

675 reads from the plasmid preparations and Nanopore reads from the plasmid preparations. B) Six

676 examples of candidate plasmids are visualized in plasmid maps. The outermost black circle indicates

677 the plasmid chromosome, the coding sequence regions are colored according to their predicted

678 function: replication (blue), mobilization (violet), transposition of DNA (green), antimicrobial

679 resistance (red), toxin-antitoxin systems (orange), hypothetical proteins (hp) and other proteins

680 (grey). The blue and green line indicate the GC and AT-content, respectively. The plasmids are

681 named according to their origin, CIV (Côte d'Ivoire), POL (Poland), USA.1 (USA), BRA (Brasil),

682 CZE (Czechia), and IND (India). Some sequencing errors might still be present in the candidate

683 plasmid sequences, which are likely the reason why a few open reading frames are not properly

684 predicted and appear fragmented, such as the gene encoding for AmpC and Macrolide efflux pump

685 genes in the plasmid from Czechia. A detailed description about the plasmids is available from

686 Figshare at https://doi.org/10.6084/m9.figshare.14039390.

687 **Table S1**. Sewage sample information.
688

689 **Figure S1: Length of nanopore sequencing reads.** The violin plot displays log transformed read

690 lengths. The horizontal dashed lines indicate log values for 1.000 and 10.000 bases length,

691 respectively. Most reads exhibit a read length below 10.000 bases, which is the cut-off value for our

692 assembly workflow, and most of the reads are between 1.000 and 10.000 bases long.

693

694 **Figure S2. Plasmid and virus (phage)-related circular DNA elements.** The bar plots display the

695 fraction (A) and total counts (B) of circular contigs containing Pfam IDs specific for plasmid and

696 virus/phage -related proteins per sample. Each predicted protein by prodigal was searched against the

697 pfam databases using HMMER hmmscan and filtered for a p-value less than 0.00001. In a small

698 subset of assemblies we identify both viral and plasmid associated genes. Pfam ID's classified as

699 "other than plasmid & viral" might still be plasmid relevant; they are just not specified as plasmid-

700 related based on the stringent scheme used.

701

702 **Figure S3. Comparison to known plasmids in plasmid database (PLSDB) – clustering on**

703 **individual plasmid level.** Samples with less than 100 circular assembled contigs where remove from

704 the analysis as well as plasmids with less than 10 occurrences over all samples. Clustering of samples

705 (columns) was done using Euclidean distance of the centered log ratio (clr)-transformed values.

25

706

707 **Figure S4. Comparison between plasmidome samples – MASH distances.** All plasmid candidate

708 sequences for each sample from the five examined continents were sketched using MASH, distances

709 calculated, and visualized by principal component analysis. A) This plot displays the differences and

710 similarities between all 24 plasmidome samples. B) This plot displays the differences and similarities

711 between 22 plasmidome samples (all samples, except NGA and BRA.2).

712

713 **Figure S5. Heatmaps depicting antimicrobial resistance profiles from the whole community**

714 **and plasmidomes from global sewage** based on presence/absence (A) and centered log ratio (clr)-

715 transformed abundancies (B) of antimicrobial resistance gene classes. The antimicrobial resistance

716 genes were identified in a ResFinder-based analysis using the Illumina reads from the whole

717 community, Illumina reads from the plasmid preparations, and Nanopore reads from the plasmid

718 preparations.

719

720 **Figure S6. Functional characterization of circular DNA elements – GO annotation.** The heat

721 map displays centered log ratio (clr)-transformed abundancies of GO annotations assigned to

722 predicted proteins. Samples with less than 100 circular assembled contig were remove from the

723 analysis as well as GO identifiers with less than 10 occurrences over all samples. The clustering of

724 samples was performed using Euclidean distance of the clr-transformed values resulting in 2 main
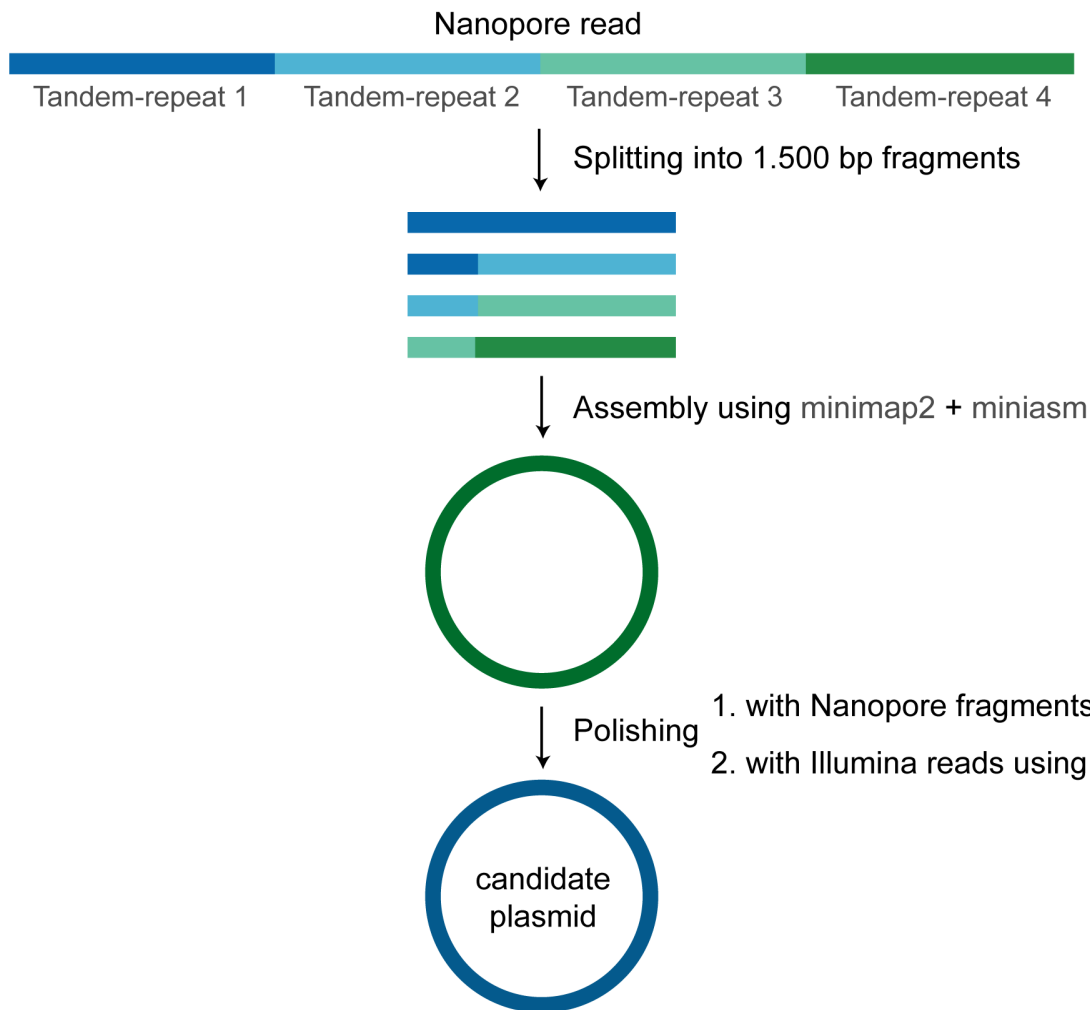
725 clusters.

726 **Figure S7. Comparison of AMR genes with prevalence data by CARD**

727 **(https://card.mcmaster.ca).** The most frequently observed AMR genes that were more abundant in

728 plasmidomes (as compared to in the whole community sequencing data) were explored at the CARD

729 website. Here, the prevalence for AMR genes is presented for a selection of pathogens, whether they

730 are associated with the plasmid or chromosome. The prevalence data are calculated as follows:

731 Antimicrobial resistance (AMR) molecular prevalence data were generated using the Resistance

732 Gene Identifier (RGI), a tool for putative AMR gene detection from submitted sequence data using

733 the AMR detection models available in CARD. To generate prevalence data, RGI was used to

734 analyze molecular sequence data available in NCBI Genomes for 88 pathogens of interest. For each

735 of these pathogens, complete chromosome sequences, complete plasmid sequences, and whole

736 genome shotgun (WGS) assemblies were analyzed individually by RGI. RGI results were then

737 aggregated to calculate percent occurrence. (See also Alcock *et al*., NAR, 2020,

738 https://academic.oup.com/nar/article/48/D1/D517/5608993).

739

740

# A



Nanopore read

Tandem-repeat 1 | Tandem-repeat 2 | Tandem-repeat 3 | Tandem-repeat 4

Splitting into 1.500 bp fragments

Assembly using minimap2 + miniasm

Polishing  1. with Nanopore fragments using racon
            2. with Illumina reads using pilon

candidate plasmid

# B



number of contigs

log2(contig length)

circular

linear

A

B

Rep_3 (PF01051)

Asia
Africa
South America
Europe
North America
Reference PF01051 domains

Aminotran_5 (PF00266.20)
AMP_N (PF05195.17)
Chlamy_scaf (PF09675.11)
DEDD_Tnp_IS110 (PF01548.18)
HTH_17 (PF12728.8)
HTH_23 (PF13384.7)
HTH_36 (PF13730.7)
HTH_Crp_2 (PF13545.7)
Mob_Pre (PF01076.20)
MobA_MobL (PF03389.16)
MobC (PF05713.12)
PilM_2 (PF11104.9)
PriCT_1 (PF08708.12)
Relaxase (PF03432.15)
Rep_1 (PF01446.18)
Rep_2 (PF01719.18)
Rep_3 (PF01051.22)
Rep_trans (PF02486.20)
RepL (PF05732.12)
Replicase (PF03090.18)
Resolvase (PF00239.22)

Toprim_2 (PF13155.7)
Transposase_20 (PF02371.17)
Gemini_AL1 (PF00799.21)
Gemini_coat (PF00844.19)
Phage_CRI (PF05144.15)
Phage_F (PF02305.18)
Phage_X (PF05155.16)
RNA_helicase (PF00910.23)
TNV_CP (PF03898.14)
Viral_Rep (PF02407.17)
other

**A** PLSD assignment

**B**