1 **Preterm infants harbour diverse *Klebsiella* populations, including atypical species that encode**

2 **and produce an array of antimicrobial resistance- and virulence-associated factors**

3

4 Yuhao Chen[1]*, Thomas C. Brook[2]*, Cho Zin Soe[3], Ian O'Neill[3], Cristina Alcon-Giner[3], Onnicha

5 Leelastwattanagul[4], Sarah Phillips[3], Shabhonam Caim[3], Paul Clarke[5,6], Lindsay J. Hall[3]†, Lesley

6 Hoyles[1,7]†

7

8 [1]Department of Surgery & Cancer, Faculty of Medicine, Imperial College London, London, United

9 Kingdom

10 [2]Department of Biomedical Sciences, Faculty of Science and Technology, University of Westminster,

11 London, United Kingdom

12 [3]Gut Microbes & Health, Quadram Institute Bioscience, Norwich Research Park, Norwich, United

13 Kingdom

14 [4]Bioinformatics & Systems Biology Program, School of Bioresources and Technology, King

15 Mongkut's University of Technology Thonburi (Bang Khun Thian Campus), Bangkok, Thailand

16 [5]Neonatal Intensive Care Unit, Norfolk and Norwich University Hospitals NHS Foundation Trust,

17 Norwich, United Kingdom

18 [6]Norwich Medical School, University of East Anglia, Norwich, United Kingdom

19 [7]Department of Biosciences, School of Science and Technology, Nottingham Trent University,

20 Nottingham, United Kingdom

21 *These authors contributed equally to this work.

22 †**Correspondence:** Lindsay J. Hall, Lindsay.Hall@quadram.ac.uk; Lesley Hoyles,

23 lesley.hoyles@ntu.ac.uk

24

25

**ABSTRACT**

*Klebsiella* spp. are frequently enriched in the gut microbiota of preterm neonates, and overgrowth is associated with necrotizing enterocolitis (NEC), nosocomial infections and late-onset sepsis. Little is known about the genomic and phenotypic characteristics of preterm-associated *Klebsiella* as previous studies have focussed on recovery of antimicrobial-resistant isolates or culture-independent molecular analyses. The aim of this study was to better characterize preterm-associated *Klebsiella* populations using phenotypic and genotypic approaches. Faecal samples from a UK cohort of healthy and sick preterm neonates ($n$=109) were screened on MacConkey agar to isolate lactose-positive *Enterobacteriaceae*. Whole-genome sequences were generated for *Klebsiella* spp., and virulence and antimicrobial resistance genes identified. Antibiotic susceptibility profiling, and *in vitro* macrophage and iron assays were undertaken for the *Klebsiella* strains. Metapangenome analyses with a manually curated genome dataset were undertaken to examine diversity of *Klebsiella oxytoca* and related bacteria in a publicly available shotgun metagenome dataset. Approximately one-tenth of faecal samples harboured *Klebsiella* spp. (*Klebsiella pneumoniae*, 7.3 %; *Klebsiella quasipneumoniae*, 0.9 %; *Klebsiella grimontii*, 2.8 %; *Klebsiella michiganensis*, 1.8 %). Isolates recovered from NEC- and sepsis-affected infants and those showing no signs of clinical infection (i.e. 'healthy') encoded multiple β-lactamases. No difference was observed between isolates recovered from 'healthy' and sick infants with respect to *in vitro* siderophore production (all encoded enterobactin in their genomes). All *K. pneumoniae*, *K. quasipneumoniae*, *K. grimontii* and *K. michiganensis* faecal isolates tested were able to reside and persist in macrophages, indicating their immune evasion abilities. Metapangenome analyses of published metagenomic data confirmed our findings regarding the presence of *K. michiganensis* in the preterm gut. There is little difference in the phenotypic and genomic characteristics of *Klebsiella* isolates recovered from 'healthy' and sick infants. Identification of β-lactamases in all isolates may prove problematic when defining treatment regimens for NEC or sepsis, and suggests 'healthy' preterm infants contribute to the resistome. Refined analyses with curated sequence databases are required when studying closely related species present in metagenomic data.

53 **Keywords:** *Klebsiella oxytoca*, shotgun metagenomics, taxonomy, microbiome.

54

55 **Author notes:** All supporting data have been provided within the article or through supplementary

56 data files, available from figshare.

57

58 **Abbreviations:** AMR, antimicrobial resistance; ANI, average nucleotide identity; CARD,

59 Comprehensive Antibiotic Resistance Database; GI, gastrointestinal; LOS, late-onset sepsis; LPE,

60 lactose-positive *Enterobacteriaceae*; LPS, lipopolysaccharide; MAG, metagenome-assembled

61 genome; NEC, necrotizing enterocolitis; NICU, neonatal intensive care unit; NNUH, Norfolk and

62 Norwich University Hospital; OTU, operational taxonomic unit; VFDB, Virulence Factors of

63 Pathogenic Bacteria Database.

64

65 **Impact statement**

66 Polyphasic characterization of isolates recovered from the faeces of preterm infants has demonstrated

67 that *Klebsiella* spp. recovered from these patients are genomically more diverse than previously

68 recognized. All *K. pneumoniae*, *K. quasipneumoniae*, *K. grimontii* and *K. michiganensis* faecal

69 isolates studied were able to reside and persist in macrophages, indicating their immune evasion

70 abilities and potential for causing infections in at-risk infants. The identification of *K. michiganensis*

71 in samples, and the abundance of *K. michiganensis* genomes in public repositories, adds to the

72 growing body of evidence indicating that *K. michiganensis* is likely to be more clinically relevant than

73 *K. oxytoca* in human-associated infections. Metapangenome analyses of publicly available shotgun

74 metagenomic data confirmed the prevalence of *K. michiganensis* in the faeces of preterm infants, and

75 highlighted the need for refined taxonomic analyses when splitting closely related species from one

76 another in metagenomic studies.

77

78 **Data summary**

79 16S rRNA gene sequence data associated with this article have been deposited at

80 DDBJ/ENA/GenBank under BioProject accession PRJEB34372. The Whole Genome Shotgun project

81    has been deposited at DDBJ/ENA/GenBank under BioProject accession PRJNA471164. The

82    metagenome-assembled genomes are available from figshare. Published sequence data of Ward *et al.*

83    (1), used to generate the metagenome-assembled genomes, are available under BioProject accession

84    number 63661.

85

86

## INTRODUCTION

The gut microbiota encompasses bacteria, archaea, lower eukaryotes and viruses, with these microbes contributing to host gastrointestinal (GI) and systemic health. Host–microbiome interactions within the intestine are particularly important in neonates, contributing to development of the immune response, establishment of the gut microbiome and protection from infections (2,3). Term infants (i.e. gestation 37 weeks) are rapidly colonised after exposure to the mother's microbiota and the environment, with streptococci and *Enterobacteriaceae* dominating in the initial phases (2), and *Bifidobacterium* spp. becoming prominent as the infant grows (2).

In contrast, colonization of preterm infants (i.e. <37 weeks' gestation) occurs in neonatal intensive care units (NICUs) and is shaped by the significant number of antibiotics ('covering' (i.e. to cover possible early onset infection from birth) and treatment) these infants receive in the first days and weeks post birth. The microbiota in preterm infants is enriched for bacteria such as *Enterobacteriaceae*, *Enterococcus* and *Staphylococcus* (4,5).

Critically, colonization of these at-risk infants with potentially pathogenic taxa, in concert with an unstable microbiome, and immaturity of their GI tract and immune system, is thought to contribute to nosocomial infections such as late-onset sepsis (LOS) or necrotizing enterocolitis (NEC) (6–12).

The family *Enterobacteriaceae* comprises more than 25 genera of catalase-positive, oxidase-negative Gram-negative bacteria and encompasses many pathogens [e.g. *Escherichia* (*Esc.*) *coli*, *Klebsiella pneumoniae*, *Shigella* (*Shi.*) *dysenteriae*, *Enterobacter* (*Ent.*) *cloacae*, *Serratia* (*Ser.*) *marcescens* and *Citrobacter* spp.] (13). While coagulase-negative staphylococci are the most common cause of LOS in preterm infants, *Enterobacteriaceae* that translocate from the preterm gut to the bloodstream also cause this condition (8,9). In addition, *Enterobacteriaceae* are associated with higher morbidity than the staphylococci, and blooms in *Proteobacteria* – thought to be linked to impaired mucosal barrier integrity – have been reported immediately prior to the diagnosis of LOS (8,9,14). Predictions made from shotgun metagenomic data show replication rates of all bacteria – and especially the *Enterobacteriaceae* and *Klebsiella* – are significantly increased immediately prior to NEC diagnosis (15). This altered gut microbiome influences intestinal homeostasis and contributes to NEC (16), in tandem with the immature preterm immune system contributing to intestinal pathology

115    in response to blooms of *Proteobacteria*.

116        Associations between *Klebsiella*-related operational taxonomic units (OTUs) and the

117    development of NEC have been noted, suggesting members of this genus contribute to the aetiology

118    of NEC in a subset of patients (17,18). Although Sim *et al.* (17) found one of their two distinct groups

119    of NEC infants had an overabundance of a *Klebsiella* OTU, these researchers failed to identify a

120    single predominant species of *Klebsiella*, recovering representatives of several genera (*K.*

121    *pneumoniae*, *Klebsiella oxytoca*, *Klebsiella aerogenes*, *Ent. cloacae*, *Esc. coli* and *Ser. marcescens*)

122    from samples. *Klebsiella* spp. and their fimbriae-encoding genes were significantly enriched in faeces

123    collected immediately prior to the onset of NEC in a US infant cohort. These fimbriae may contribute

124    to the overexpression of TLR4 receptors observed in preterm infants (15). Confirming the role of

125    these bacteria in NEC will require reproducing certain aspects of the disease in model systems, using

126    well-characterized bacteria recovered from preterm infants (14,17).

127        To date, there is limited information on the genomic and phenotypic features of preterm-

128    associated *Klebsiella* spp. Thus, to characterise these important opportunistic pathogens, and to build

129    a collection of preterm-associated *Klebsiella* strains for use in future mechanistic studies relevant to

130    preterm-infant health, we isolated and characterized (phenotypically and genomically) bacteria from a

131    cohort of preterm neonates enrolled in a study at the Norfolk and Norwich University Hospital

132    (NNUH), Norwich, United Kingdom. Recovered *Klebsiella* isolates were subject to additional

133    phenotypic tests that complemented genomic data. In addition, for the increasingly important species

134    *K. oxytoca*, in which our laboratories have a specific interest, metapangenome analyses were

135    undertaken to better understand the prevalence and potential virulence of this organism and related

136    species in the context of the preterm neonate gut microbiota.

137

138    **METHODS**

139    **Collection of faecal samples.** Faeces were collected from premature neonates (<37 weeks' gestation)

140    (**Supplementary Table 1**). The Ethics Committee of the Faculty of Medical and Health Sciences of

141    the University of East Anglia (Norwich, UK) approved this study. The protocol for faeces collection

142    was laid out by the Norwich Research Park (NRP) Biorepository (Norwich, UK) and was in

143  accordance with the terms of the Human Tissue Act 2004 (HTA), and approved with licence number

144  11208 by the Human Tissue Authority. Infants admitted to the NICU of the NNUH were recruited by

145  doctors or nurses with informed and written consent obtained from parents. Collection of faecal

146  samples was carried out by clinical researchers and/or research nurses, with samples stored at -80 °C

147  prior to DNA extraction.

148

149  **16S rRNA gene sequencing and analyses.** DNA was extracted from samples using the FastDNA

150  SPIN Kit for Soil (MP Biomedicals) and processed for sequencing and analyses as described

151  previously (19). This 16S rRNA gene sequence data associated with this project have been deposited

152  at DDBJ/ENA/GenBank under BioProject accession PRJEB34372.

153

154  **Isolation of bacteria and biochemical characterization.** For isolation work, a single faecal sample

155  from each baby ($n$=109; **Supplementary Table 1**) was thawed and 0.1 g homogenised in 1 mL TBT

156  buffer (100 mM Tris/HCl, pH 8.0; 100 mM NaCl; 10 mM $MgCl_2 \cdot 6H_2O$). Homogenates were serially

157  diluted $10^{-1}$ to $10^{-4}$ in TBT buffer. Aliquots (50 µL) of homogenate were spread on MacConkey agar

158  no. 3 (Oxoid Ltd) plates in triplicate and incubated aerobically at 37 °C overnight.

159  Differential counts (based on colony morphology) of all lactose-positive (i.e. pink) colonies

160  were made in triplicate to calculate colony-forming units (CFUs) per gramme wet-weight faeces. One

161  of each colony type per plate was selected and re-streaked on MacConkey agar three times to purify,

162  incubating aerobically at 37 °C overnight each time. A single colony from each pure culture was re-

163  suspended in 5 mL of sterile distilled water; the API 20E kit (bioMérieux) was used according to the

164  manufacturer's instructions to give preliminary identities for each of the isolates recovered.

165

166  **DNA extraction, whole-genome sequencing and assembly.** DNA was extracted using a phenol–

167  chloroform method fully described previously (20) from overnight cultures of strains, and sequenced

168  using the 96-plex Illumina HiSeq 2500 platform to generate 125 bp paired-end reads (21). Raw data

169  provided by the sequencing centre were checked using fastqc v0.11.4

170  (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/); no adapter trimming was required, and

171  reads had an average Phred score >25. MetaPhlAn2.6 (22) was used to identify which species genome

172  sequences represented. According to the results given by MetaPhlAn2.6, appropriate reference

173  genomes were retrieved from Ensembl Genome (http://bacteria.ensembl.org/index.html) to guide

174  reference-based assembly using BugBuilder v1.0.3b1 (default settings for Illumina data) (23).

175  Summary statistics for the *Klebsiella* genome sequences generated in this study, including accession

176  numbers, can be found in **Supplementary Table 2**. This Whole Genome Shotgun project has been

177  deposited at DDBJ/ENA/GenBank under BioProject accession PRJNA471164.

178

179  **Genome analyses.** Average nucleotide identity (ANI) between genome sequences of isolates and

180  reference strains (*Klebsiella grimontii* 06D021$^T$, GCA_900200035; *K. oxytoca* 2880STDY5682490,

181  GCA_900083895.1; *Klebsiella michiganensis* DSM 25444$^T$, GCA_002925905) was determined using

182  FastANI (default settings) (24).

183      *K. oxytoca*, *K. michiganensis* and *K. grimontii* genomes were uploaded to the *Klebsiella*

184  *oxytoca*/*michiganensis* MLST website (https://pubmlst.org/koxytoca/) sited at the University of

185  Oxford (25) on 28 July 2019 to determine allele number against previously defined house-keeping

186  genes (*rpoB*, *gapA*, *mdh*, *pgi*, *phoE*, *infB* and *tonB*). *K. pneumoniae* genomes were analysed using the

187  Institut Pasteur MLST database (https://bigsdb.pasteur.fr/klebsiella/klebsiella.html). Kleborate (26,27)

188  and Kaptive (28) were used to identify capsular type and O antigen type.

189      Virulence genes were identified by BLASTP of genome amino acid sequences against the

190  Virulence Factors of Pathogenic Bacteria Database (VFDB; 'core dataset' downloaded 27 July 2019)

191  (29); results are reported for >70 % identity and 90 % query coverage. Antimicrobial resistance

192  (AMR) genes were identified by BLASTP against the Comprehensive Antibiotic Resistance Database

193  (CARD) download (27 July 2019; protein homolog dataset) (30); only strict and perfect matches with

194  respect to CARD database coverage and bit-score cut-off recommendations are reported.

195      Genomic traits were visualized using anvi'o-5.5 (31) according to the pangenomic workflow.

196  Briefly, for each figure presented herein, genomes were used to create an anvi'o contigs database,

197  which contained ORFs predicted using Prodigal v2.6.3 (32). A multiple-sequence alignment was

198  created using BLASTP. Markov CL algorithm (33) was used to identify gene clusters (--mcl-inflation

199    10; high sensitivity for identifying gene clusters of closely related species or strain level). Gene

200    clusters and genomes were organized using Euclidean distance and Ward linkage, with results

201    visualized using GoogleChrome.

202

203    **Phenotypic characterization of *Klebsiella* isolates**

204    <u>**Iron assays.**</u> Pre-cultures of *Klebsiella* isolates (5 mL) were grown overnight in LB broth (37 °C, 160

205    rpm). Aliquots (500 µL) were harvested (4000 rpm, 20 min) and the cell pellets washed twice with

206    PBS. The cell suspensions (50 µL) were used to inoculate 5 mL cultures containing M9 minimal

207    medium ($Na_2HPO_4$, 6.9 g/L; $KH_2PO_4$, 3 g/L; NaCl, 0.5 g/L; $NH_4Cl$, 1 g/L; $CaCl_2$, 0.1 mM; $MgSO_4$, 2

208    mM; 0.2 % glucose) at 37 °C. At 20 h, bacterial growth and siderophore production were measured

209    using the CAS assay (34). An aliquot (100 µL) of the cell culture supernatant was mixed with CAS

210    dye (100 µL), followed by the shuttle solution (4 µL) and siderophore production monitored at 620

211    nm at 4 h using a BioRad Benchmark Plus microplate spectrophotometer. A decrease in the blue

212    colour of the CAS dye was measured using uninoculated medium as control. The estimated amount of

213    total siderophore produced by *Klebsiella* isolates was calculated using the CAS standard curve based

214    upon a desferrioxamine B standard (1:1).

215

216    <u>**Macrophage assays.**</u> All strains were grown on LB broth + 1.5 % agar and incubated overnight at 37

217    °C. THP-1 monocytes were obtained from ATCC (TIB-202) and were maintained in RPMI (Gibco:

218    72400021) plus 10 % heat-inactivated foetal bovine serum (FBS; Gibco: 10500064) in a humidified

219    incubator at 37 °C with 5 % $CO_2$. THP-1 monocytes were differentiated into macrophages in RPMI +

220    10 mM HEPES + 10 % FBS + 10 ng/mL phorbol 12-myristate 13-acetate (PMA; Sigma) and seeded

221    at $1 \times 10^5$ cells per well of a 96-well tissue culture dish and incubated for 15 h. Overnight cultures of

222    bacteria were diluted 1:100 into fresh LB broth and grown until mid-exponential phase. Bacteria were

223    then washed twice with PBS and diluted to $1 \times 10^7$ cfu/mL in RPMI + 10 mM HEPES + 10 % FBS and

224    100 µL of bacteria was added to each well. Plates were then centrifuged at 300 *g* for 5 min to

225    synchronize infections. Bacteria–macrophage co-culture was incubated at 37 °C/5 % $CO_2$ for 30 min

226    to allow for phagocytosis. Cells were then washed three times in PBS and medium was replaced with

227     above culture medium supplemented with 300 μg/mL gentamicin and 100 units/mL polymyxin B to

228     eliminate extracellular bacteria. Cell were again incubated at 37 °C/5 % $CO_2$ for 1.5 h. Cells were then

229     washed three times with PBS and medium for cells for later time points was replaced with culture

230     medium supplemented with 300 μg/mL gentamicin and incubated for a further 4.5 h. Intracellular

231     bacterial load was enumerated by lysing macrophages in PBS + 1% Triton X-100 for 10 min at room

232     temperature, serially diluting cultures and plating on LB agar. Plates were incubated overnight at 37

233     °C and colonies counted the following day.

234

235     **Calculation of antibiotic minimal inhibitory concentration (MIC) for the *Klebsiella* isolates.**

236     Broth microdilution method was used to calculate the MIC of the *Klebsiella* isolates. Serial two-fold

237     dilutions of benzylpenicillin, gentamicin, and meropenem were added to sterile nutrient broth. The

238     antibiotics used in this assay were supplied by the NICU of NNUH. Inoculum for each of the isolates

239     was prepared using 10 mL from a fresh overnight culture. Microplates were incubated for 24 h at 37

240     °C under aerobic conditions. Optical density was monitored using a plate reader (BMG Labtech, UK)

241     at 595 nm. MICs were determined as the lowest concentration of antibiotic inhibiting any bacterial

242     growth. All experiments were repeated in triplicate. For the aminoglycoside gentamicin and the

243     carbapenem meropenem, *Klebsiella* (*Enterobacteriaceae*) breakpoints were determined according to

244     European Committee on Antimicrobial Susceptibility Testing (EUCAST) guidelines (version 8.1,

245     published 16 May 2018,

246     http://www.eucast.org/fileadmin/src/media/PDFs/EUCAST_files/Breakpoint_tables/v_8.1_Breakpoin

247     t_Tables.pdf). No EUCAST data were available for benzylpenicillin (EUCAST states this

248     aminopenicillin has no clinically useful activity against *Enterobacteriaceae*).

249

250     **Estimation of abundance of *K. oxytoca* in shotgun metagenomic data.** We chose to analyse a

251     published preterm gut metagenome dataset (1) in this study, as it had been previously used to identify

252     associations between uropathogenic *Esc. coli* and NEC. Trimmed, human-filtered, paired-end read

253     data deposited in the Sequence Read Archive by Ward *et al.* (1) are available under BioProject

254     accession number 63661. Information on Ward samples included in this study can be found in

255    **Supplementary Table 3**. Ward *et al.* (1) used MetaPhlAn to determine abundance of bacteria in

256    samples. However, the marker genes used to enumerate *K. oxytoca* in the MetaPhlAn2.6 database are

257    derived from 11 genomes, five of which are not *K. oxytoca* (*Raoultella ornithinolytica* 10-5246,

258    GCF_000247895; *K. michiganensis* E718, GCF_000276705; *K. michiganensis* Kleb_oxyt_10-

259    5250_V1, GCF_000247915; *K. michiganensis* KCTC 1686, GCF_000240325; *K. michiganensis*

260    Kleb_oxyt_10-5242_V1, GCF_000247835). Therefore, relative abundance of bacteria was instead

261    determined using Centrifuge (35). While MetaPhlAn2.6 relies on a pre-compiled database of unique

262    marker genes for determining taxonomic abundance, the Centrifuge database can be updated at will

263    using genomes downloaded from NCBI. A bacteria- and archaea-specific complete genome database

264    was generated for use with Centrifuge via NCBI on 1 July 2018. Species-level abundances, based on

265    read-level data, for *K. oxytoca* and *K. michiganensis* in the study of Ward *et al.* (1) were determined.

266    (NB: *K. grimontii* genomes were not included in the Centrifuge database, nor are they included in

267    MetaPhlAn2.6 or the most-recent version of Kraken2.)

268

269    **Metapangenome analyses of *K. oxytoca*, *K. michiganensis* and *K. grimontii*.** A total of 162 *K.*

270    *oxytoca*-related whole-genome sequences were retrieved from GenBank on 31 May 2018

271    (**Supplementary Table 4**). On the basis of *bla*$_{OXY}$, phylogenetic and ANI analyses (36), these had

272    been confirmed to belong to *K. oxytoca* (*n*=64), *K. grimontii* (*n*=24) and *K. michiganensis* (*n*=74)

273    (**Supplementary Table 5**). Prokka v1.13.3 (37) was used to annotate the 162 downloaded and five

274    infant genomes. The resulting .gff files were subject to pangenome analyses using Roary v3.12.0

275    (default settings) (38). Genes present in 165–167 strains were defined as the core cluster, while those

276    present in 25–164 strains were defined as the accessory cluster. Remaining genes that only existed in

277    single strains were classified into strain-specific clusters. FastTree v2.1.10 (39) was used to generate a

278    phylogenetic tree from the core gene alignment, with the tree visualized using FigTree v1.4.4

279    (http://tree.bio.ed.ac.uk/software/figtree/).

280          PanPhlAn (panphlan_pangenome_generation.py v1.2.3.6; panphlan_profile.py v1.2.2.3;

281    panphlan_map.py v1.2.2.5; (40)) was used to profile strains within metagenomes using the Roary-

282    generated pangenome dataset. Gene-family clusters across all 167 available genomes and centroid

283    sequence files outputted from Roary were uploaded to PanPhlAn to build a Bowtie2-indexed

284    pangenome database, against which raw reads (concatenated read pair files) were mapped using

285    Bowtie2 v2.3.0. The coverages of all gene positions were detected and extracted using samtools

286    v1.4.1, and then integrated to a gene-family coverage profile for each sample. Pangenome-related

287    strains were predicted to exist if a consistent and similar coverage depth across a set of 5774 gene

288    families was detected under the non-default parameters of PanPhlAn (--min_coverage 1, --left_max

289    1.70, --right_min 0.30; panphlan_profile.py). Principal component analysis (PCA) was performed on

290    500 accessory genes randomly selected from the pangenome using the R package FactoMineR (41),

291    allowing us to distinguish different species at the gene level.

292

293    **Recovery of metagenome-assembled genomes (MAGs) from metagenomes.** For metagenome

294    samples in which *K. oxytoca*-related strains were identified using PanPhlAn, attempts were made to

295    recover them as MAGs. All reads in samples were mapped against the pangenome database using

296    Bowtie2 and mapped paired-end reads were extracted by using FastQ Screen v0.11.3 as new fastq

297    files (with parameter --tag, --filter 3). The extracted paired-end reads were assembled using SPAdes

298    v3.12.0 (42). These assemblies were known as original MAGs. A genome size of 5.5 Mb was set as a

299    strict threshold: any assembly whose genome size was lower than this threshold was not considered in

300    downstream analyses. FastANI was applied to calculate the ANI (cut-off 95 %) between MAGs and

301    the three species reference genomes to double-check the predominant species in corresponding

302    samples. The quality of each MAG was assessed using CheckM v1.0.18 (43).

303         A small number of the original MAGs were of high quality (44), but some contained a large

304    number of contaminant contigs. All original MAGs were decontaminated as follows. Coding

305    sequences of species-specific genomes in the pangenome were predicted using Prodigal v2.6.3 (Hyatt

306    et al., 2010) with default settings and resulting multi-FASTA files containing protein sequences were

307    concatenated to single files, which were used to build *K. oxytoca*-, *K. michiganensis*- and *K.*

308    *grimontii*-specific databases in Diamond (45) format. Contigs of the original MAGs were aligned

309    against the corresponding database under different minimum identity (%) cutoffs to report sequence

310    alignments (with parameter --id 95, 96, 97, 98, 99 and 100). All unmapped contigs and contigs <500

311    nt in length were discarded from the original MAGs and the quality of new MAGs acquired was

312    evaluated again using CheckM, to identify a Diamond BLAST identity threshold at which

313    decontamination was effective while maintaining high genome completeness.

314

315    **RESULTS AND DISCUSSION**

316    **Composition of the microbiota of preterm neonates**

317        First faecal samples available after birth were collected from 109 hospitalized preterm infants

318    ($n = 50$ female; $n = 59$ male) in the NICU of the NNUH (**Supplementary Table 1**; **Figure 1a**). On

319    the basis of 16S rRNA gene sequencing data, *Enterobacteriaceae* were detected in the faeces of 42

320    (38.5 %) of the infants (**Figure 1b**).

321        All faecal samples were screened for *Enterobacteriaceae* using MacConkey agar no. 3

322    (**Figure 1c**). Forty-six (42.2 %) samples were positive for *Enterobacteriaceae* ($n = 9$ lactose-negative,

323    carriage rate 8.3 %; $n = 37$ lactose-positive, carriage rate 33.9 %). Lactose-negative isolates were not

324    characterized further. API 20E was used to provide tentative identification of LPE from 36 neonates

325    (isolate could not be resuscitated for neonate P054) (**Figure 1d**). Of the 36 infants from whose faeces

326    isolates were recovered, 23 were healthy, three had or were subsequently diagnosed with NEC, eight

327    had suspected sepsis, one had an operation for gastroschisis and one was diagnosed with an eye

328    infection after the faecal sample was taken (**Supplementary Table 1**).

329

330    **Whole-genome sequencing of neonatal faecal LPE**

331        Whole-genome sequences were obtained for 56 LPE. MetaPhlAn2.6 was used to assign

332    identities to genomes (not shown). Among the 56 isolates sequenced, 20 were identified as *K.*

333    *pneumoniae* (carriage rate 8.3 %; **Supplementary Table 2**), 14 were *Ent. cloacae* complex (carriage

334    rate 11.9 %), 13 were *Esc. coli* (carriage rate 11.9 %), five were *K. oxytoca* (**Supplementary Table**

335    **2**), two were *Citrobacter freundii* (carriage rate 1.8 %; (46)), one was *Citrobacter murliniae* (carriage

336    rate 0.9 %; (46)) and one was *R. ornithinolytica* (carriage rate 0.9 %; (47)). *Esc. coli* and *Ent. cloacae*

337    complex isolates will be discussed in detail elsewhere. MetaPhlAn2.6-generated identities matched

338    those given by API 20E (**Supplementary Table 2**).

339        Reference-based assembly of genomes was performed using BugBuilder (23)

340    (**Supplementary Table 2**). To determine whether preterm neonates may harbour more than one strain

341    of a species in their faecal microbiota, nine isolates (#64–#73) were collected from neonate P008.

342    These had all been identified as *K. pneumoniae* by API 20E and genome data. ANI across the nine

343    isolates was >99.99 %. To determine whether the isolates were identical, gene content analysis was

344    performed using Roary (38). The average number of CDSs among these isolates was 5,385 (6.25)

345    (**Supplementary Figure 1a–d**). Anvi'o showed the isolates were highly similar (**Supplementary**

346    **Figure 1e**). Isolates of the same species from other neonates were also found to be identical to one

347    another (#102 and #103 from P080; #118 and #119 from P124). For sets of identical isolates, only one

348    was taken forward for further analyses. This left 14 distinct *Klebsiella* strains (9 *K. pneumoniae*; 5 *K.*

349    *oxytoca*) for further analyses.

350

351    **Genome analyses of *K. pneumoniae* strains**

352        *K. pneumoniae* is a commensal of the human gut microbiota and can cause nosocomial

353    infections, NEC and LOS in premature neonates (9,17,48–50),(51). The genetic backgrounds of the

354    neonate isolates were explored, to determine virulence and AMR genes encoded within the strains'

355    genomes.

356        Each isolate was genetically different: i.e. no two infants harboured the same strain of *K.*

357    *pneumoniae* (**Supplementary Figure 2**). ANI analyses with representative strains of the seven

358    phylogenetic groups of *K. pneumoniae* (52) showed eight of the neonatal isolates were *K. pneumoniae*

359    (98.83–98.98 % ANI with *K. pneumoniae* ATCC 13883$^T$ (GCA_000742135)) and one (#91) was *K.*

360    *quasipneumoniae* (98.5 % ANI with *Klebsiella quasipneumoniae* subsp. *quasipneumoniae* 01A030$^T$

361    (GCA_000751755)). MLST identified six STs within the *K. pneumoniae* strains (**Figure 2a**). *K.*

362    *quasipneumoniae* #91 had a novel *mdh* allele, so no ST could be specified for this strain. None of the

363    STs belonged to clonal complex (CC) 258, responsible for hospital outbreaks due to its frequent

364    carriage of KPC and other acquired AMR genes (53). The capsule of *K. pneumoniae* and related

365    species is considered one of its major virulence factors. K1, K2 and K5 capsular types and

366    hypervirulent types have strong associations with human infectious diseases (54,55). None of our

367    neonatal isolates had a capsular type commonly associated with infections or hypervirulent *K.*

368    *pneumoniae*, though K7, K10, K11, K16 and K38 isolates have previously been recovered from

369    clinical samples in Taiwan (56). Although the capsular type of strain #74 was identified as K62 with

370    99.33 % confidence and 100 % coverage, there was one gene (*KL62-12*, according to Kaptive)

371    missing from the locus, leaving it a non-perfect match. Further analyses showed the genes associated

372    with K62 to be disrupted in strain #74 and not encoded on a contiguous stretch of DNA

373    (**Supplementary Figure 3**). Of the nine *K. pneumoniae* strains analysed using Kleborate (26,27),

374    O1v1 and O1v2 were represented equally among the O-antigen types (*n* =4 for both). These can be

375    distinguished using genomic data but are serologically cross-reactive (27). *K. quasipneumoniae* #91

376    was O3/O3a.

377        The vast majority (e.g. ~90 % in the NNUH NICU) of preterm infants receive antibiotics during

378    their NICU stay, often started routinely from admission (i.e. 'covering') if they are born very

379    premature and/or very low birth weight. Administration of antibiotics can lead to disruption of early

380    colonization by microbes, potentially encouraging growth of opportunistic pathogens such as LPE,

381    creating a selection pressure that may promote development of AMR. All strains encoded homologues

382    of *acrB*, *acrD*, *marA*, *emrR*, *CRP*, *mdtB*, *mdtC*, *baeR*, *Escherichia mdfA*, PmrF, *msbA*, OmpK37,

383    KpnE, KpnF, KpnG, KpnH and *Escherichia ampH* β-lactamase, associated with antibiotic efflux and

384    its regulation (*acrB*, *acrD*, *marA*, *emrR*, *CRP*, *mdtB*, *mdtC*, *baeR*) (57–59) and resistance to:

385    aminoglycosides; cationic antimicrobial peptides and antibiotics such as polymyxin (PmrF (60));

386    chloramphenicol (*mdfA* (61)), cefotaxime and cefoxitin (OmpK37); cefepime, ceftriaxone, colistin,

387    erythromycin, rifampin, tetracycline, streptomycin as well as enhanced sensitivity toward sodium

388    dodecyl sulfate, deoxycholate, dyes, benzalkonium chloride, chlorohexidine, and triclosan (KpnE,

389    KpnF (62)); azithromycin, ceftazidime, ciprofloxacin, ertapenem, erythromycin, gentamicin,

390    imipenem, ticarcillin, norfloxacin, polymyxin-B, piperacillin, spectinomycin, tobramycin, and

391    streptomycin (KpnG, KpnH (63)); β-lactams and penicillin (*Escherichia ampH* β-lactamase (64)). As

392    expected, core genes $bla_{SHV}$ and $bla_{OKP}$, respectively, were found in *K. pneumoniae* and *K.*

393    *quasipneumoniae* genomes (53). *K. quasipneumoniae* #91 also encoded homologues of the acquired

394    AMR gene *emrB* (a translocase that recognizes substrates including carbonyl cyanide *m*-

395    chlorophenylhydrazone, nalidixic acid and thioloactomycin) and *bacA* (confers resistance to

396    bacitracin) (**Figure 2a**). Plasmid-encoded $bla_{SHV}$ enzymes represent an important subgroup of class A

397    β-lactamases, while chromosomally encoded β-lactamase $bla_{OKP}$ cannot hydrolyse extended-spectrum

398    cephalosporins (65). Homologues of the core AMR genes *oqxA* and *oqxB* (encoding OqxAB, a

399    plasmid-encoded efflux pump that confers resistance to fluoroquinolones) were encoded by #64, #74,

400    #91, #95, #115 and #118. Strains #64 and #95 encoded homologues of FosA6, while #74, #85, #92,

401    #115 and #118 encoded homologues of FosA5 (both gene products confer resistance to fosfomycin,

402    and are core AMR genes (53)).

403        While the majority of the neonatal *K. pneumoniae* strains did not represent known pathogenic

404    lineages, virulence factors were detected in their genomes using VFDB (**Figure 2b**). The host limits

405    iron availability within the GI tract to prevent colonization by pathogens and bacterial overgrowth.

406    However, *Klebsiella* spp. have evolved numerous mechanisms to circumvent these defences. Thus,

407    we determined whether gene clusters associated with iron uptake and siderophore systems (i.e.

408    enterobactin, yersiniabactin, aerobactin, colibactin, salmochelin) were present in the strains. All

409    strains encoded enterobactin, while only #115 encoded an additional system (yersiniabactin) (**Figure

410    2b**). All strains encoded *Esc. coli* common pilus, OmpA, Hsp60, type 3 fimbriae, ferric aerobactin

411    receptor *IutA* and the AcrAB efflux pump. All strains except #102 encoded type 1 fimbriae; all strains

412    except *K. quasipneumoniae* #91 encoded typical *K. pneumoniae* lipopolysaccharide (LPS) according

413    to VFDB (>70 % amino acid identity and 90 % query coverage). *K. quasipneumoniae* #91 encoded

414    thin aggregative fimbriae, associated with biofilm formation and adhering to human mucosal or

415    epithelial surfaces. Incomplete coverage of *Klebsiella* capsule genes is likely due to the limited

416    database of VFDB compared with those used to populate Kaptive and Kleborate. Kaptive had shown

417    #91 to be K11 and O3/O3a. The core LPS region of #91 was identified using the *waa* gene cluster

418    (66); WaaL clustering with an 80 % threshold showed the strain had LPS core type 1 (67)

419    (**Supplementary Figure 4**).

420

421    **Whole-genome analyses of isolates tentatively identified as *K. oxytoca***

422    *K. oxytoca* is a minor member of the human gut microbiota, recovered at low levels from the

423    faeces of 1.6–9 % of healthy adults (68). Toxigenic *K. oxytoca* is a causative agent of antibiotic-

424    associated haemorrhagic colitis, a condition affecting mainly young and otherwise healthy outpatients

425    after brief treatment with penicillin derivatives (69). *K. oxytoca* has been detected in the faeces of a

426    subset of preterm infants via cultivation or shotgun metagenomics, but its association with preterm-

427    associated infections is unknown (1,17,70). At the DNA level, bacteria characterized phenotypically

428    as *K. oxytoca* actually represent six phylogroups/distinct species: Ko1, *K. michiganensis*; Ko2, *K.*

429    *oxytoca*; Ko3, '*Klebsiella spallanzanii*'; Ko4, '*Klebsiella pasteurii*'; Ko6, *K. grimontii*; Ko8, *K.*

430    *huaxiensis* (71–74)(75). *K. michiganensis* and *K. oxytoca* are distinguishable based on the $bla_{OXY}$ gene

431    they carry ($bla_{OXY-1}$ and $bla_{OXY-2}$, respectively) (71). *K. grimontii* was recently described to

432    accommodate Ko6 strains based on *rpoB*, *gyrA* and *rrs* gene sequences (73). All six members of the

433    complex can be differentiated by MALDI-TOF (75), but reference databases currently in routine

434    clinical use lack reference spectra of the different species to allow identification beyond *K. oxytoca*.

435    Consequently, reports on complex members other than *K. oxytoca* have only recently begun to appear

436    in the literature (73–77). The colonization of humans with *K. oxytoca* phylogroups has previously

437    been associated with the genetic backgrounds of strains: Ko2 mainly inhabits the lower GI tract, with

438    Ko1 and Ko6 generally associated with respiratory isolates and faecal isolates, respectively (72).

439    On the basis of API 20E data and initial genome (MetaPhlAn2.6) analysis, five neonatal

440    isolates were identified as *K. oxytoca* (#80, #83, #88, #99, #108). It has recently been shown that API

441    20E and MALDI-TOF using current clinical reference databases are as effective as one another for

442    characterization of complex members as *K. oxytoca* (78). MetaPhlAn2.6 cannot distinguish among

443    species of the *K. oxytoca* complex. ANI of the genomes against reference genomes showed #88 and

444    #108 to be *K. michiganensis* (both 98.78 and 98.94 % ANI, respectively with GCA_002925905) and

445    #80, #83 and #99 to be *K. grimontii* (99.18, 99.23 and 99.17 % ANI, respectively, with

446    GCA_900200035) (**Supplementary Figure 5a**), with ANI cut-off values well above the ~95 %

447    proposed for species delineation (79–81) and used by Passet & Brisse (73) to separate *K. grimontii*

448    from *K. oxytoca* and *K. michiganensis*. Phylogenetic analysis with a panel of authentic *K. oxytoca*, *K.*

449    *grimontii* and *K. michiganensis* genomes confirmed the species affiliations of the infant isolates

450    (**Supplementary Figure 5b**). Similar to the *K. pneumoniae* isolates, no two infants harboured the

451    same strain of *K. michiganensis* or *K. grimontii* (**Figure 3a**).

452    It is notable that of the publicly available genomes deposited as *K. oxytoca*, 74 were found to

453    represent *K. michiganensis*, 64 were *K. oxytoca* and 24 were *K. grimontii*. This suggests that *K.*

454    *michiganensis* may be more clinically relevant than *K. oxytoca sensu stricto*. *K. michiganensis* was

455    originally proposed to describe an isolate closely related to *K. oxytoca* recovered from a toothbrush

456    holder (74). The bacterium is now recognised as an emerging pathogen, with this recognition due to

457    improved genomic characterization of clinical isolates that would have previously been described as

458    *K. oxytoca* based on simple phenotypic tests or MALDI-TOF (78,82–84).

459

460    **Predicted virulence and AMR determinant genes of infant-associated *K. michiganensis* and *K.***

461    ***grimontii***

462    The *K. michiganensis* and *K. grimontii* strains were examined for the presence of virulence-

463    associated loci found in *K. pneumoniae* strains (53) (**Figure 3b**). Enterobactin was encoded by all

464    strains. Yersiniabactin was predicted to be encoded by *K. grimontii* #99 and *K. michiganensis* #88 and

465    #108. Other siderophore-associated gene clusters (aerobactin, colibactin and salmochelin) found in *K.*

466    *pneumoniae* were absent. An allantoinase gene cluster (including *allB/C/R/A/S* and *ybbW*), which

467    plays a role in *K. pneumoniae* liver infection (85), was identified in the three *K. grimontii* strains.

468    Due to the clinical importance of AMR in *Enterobacteriaceae*, an *in silico* AMR gene profile

469    was established for the *K. michiganensis* and *K. grimontii* strains. Homologues of 18 AMR

470    determinant genes (acquired AMR genes – *emrB*, *emrR*; core AMR genes *acrB*, *acrD*, CRP, *marA*,

471    *mdtB*, *mdtC*, *baeR*, FosA5, *pmrF*, *oqxA*, *oqxB*, *msbA*, KpnE, KpnF, KpnG, *Escherichia ampH* β-

472    lactamase) were common to the five strains, similar to the *K. pneumoniae* isolates. Both *K.*

473    *michiganensis* strains encoded homologues of OXY-1-2 (β-lactamase specific to *K. michiganensis*

474    (Ko1; (86)) and *bacA*, while #108 encoded a homologue of *aph(3')-la* (aminoglycoside

475    phosphotransferase). All *K. grimontii* strains encoded *mdtN* (potentially involved in resistance to

476    puromycin, acriflavine and tetraphenylarsonium chloride), while #83 and #99 encoded homologues of

477    OXY-6-2 (β-lactamase specific to *K. grimontii* (Ko6; (86)).

478

479    **Phenotypic characterization of *K. pneumoniae*, *K. quasipneumoniae*, *K. michiganensis* and *K.**

480    ***grimontii* neonatal isolates**

481       Five of the 13 *Klebsiella* strains we characterised were isolated from preterm infants who had

482    been diagnosed with either NEC or sepsis. Thus, we sought to link our genotypic analyses with

483    clinically important virulence traits including the ability to survive and replicate in host immune cells

484    (i.e. macrophages) and the ability to produce iron-acquiring siderophores. We also determined the

485    strains' AMR profiles for a limited set of antimicrobials.

486       Previous studies have indicated that respiratory infection-associated *K. pneumoniae* are able to

487    survive within macrophages, a critical innate immune cell type required for optimal pathogen

488    clearance (87). However, to date there is limited information relating to this ability in gut-associated

489    strains, and there is no information on other *Klebsiella* species. Thus, all *Klebsiella* strains isolated in

490    this study were tested in PMA-differentiated THP-1 macrophages using a gentamicin protection

491    assay. All strains appeared to persist within macrophages, as bacterial load was either maintained over

492    the time-course or increased or decreased between 1.5 h and 6 h, although these values were not

493    statistically significant (**Figure 4a**). These data suggest that all *Klebsiella* strains tested can reside and

494    persist in macrophages. This ability of all strains to survive, and in some cases potentially replicate,

495    within macrophages indicates their immune evasion capabilities, which may link to increased risk and

496    incidence of NEC and sepsis if these strains translocate from the 'leaky' preterm GI tract to systemic

497    sites contributing to the inflammatory cascades characteristic of these conditions.

498       Iron is a vital nutrient that performs multiple roles in cellular processes, ranging from DNA

499    replication and cell growth to protection against oxidative stress. In the healthy host, the majority of

500    iron is bound with intracellular proteins and the remaining free iron is extracellular and insoluble,

501    hence difficult to access (88). For invading pathogens, siderophore systems are critical for iron

502    competition and uptake to accomplish colonization and cause infections, this is particularly true in the

503    preterm GI tract. Preterm infants in NICU are heavily supplemented with iron as they receive many

504 red-blood-cell transfusions (increasing hepatic iron stores), iron-supplemented parenteral nutrition,

505 and supplementary oral iron within a few weeks of birth. During infection, *Klebsiella* secretes

506 siderophores to sequester iron and to establish colonization in the host. Enterobactin is the most well-

507 known siderophore produced by *K. pneumoniae* and related species, and was found to be encoded by

508 all our strains (**Figure 2b**, **Figure 3b**). The host innate immune protein lipocalin 2 binds to

509 enterobactin and disrupts bacterial iron uptake (89). *Klebsiella* species have evolved to hoodwink this

510 host response by producing several evasive siderophores (90,91). Siderophore production of

511 *Klebsiella* isolates was monitored using CAS liquid assay. All isolates tested grown in M9 minimal

512 medium were CAS-positive with the estimated siderophore concentration ranging between of 3.5 and

513 6 nM (**Figure 4b**). There was no significant difference in siderophore production between 'healthy'

514 and NEC- and sepsis-associated isolates.

515       *Klebsiella* is of concern within an AMR context, particularly in at-risk neonates, due to the

516 increasing emergence of multidrug-resistant isolates that cause severe infection (92). A UK study in

517 which 24 % of all LOS cases were caused by *Enterobacteriaceae* (8.9 % of all caused by *Klebsiella*

518 spp.) showed a high proportion (14 % and 34 %, respectively) of *Enterobacteriaceae* isolates

519 recovered from sick infants were resistant to flucloxacillin/gentamicin and amoxicillin/cefotaxime, the

520 two most commonly used empiric antibiotic combinations (51). Thus, to demonstrate antibiotic-

521 resistance phenotypes in *Klebsiella* spp. correlating to presence of AMR genotypes, we tested the

522 susceptibility of the isolates with three antibiotics commonly prescribed in NICUs; gentamicin,

523 meropenem and benzylpenicillin (**Table 1**). One strain of *K. grimontii* (#80) was potentially sensitive

524 to benzylpenicillin, an aminopenicillin currently not recognized as being clinically useful against

525 *Enterobacteriaceae*.

526       Isolates #64, #83 and #108 (all encoding KpnG and KpnH; **Figures 2a, 3b**) were resistant to

527 gentamicin, while #80, #88, #95 and #99 (all encoding KpnG and KpnH; **Figures 2a, 3b**) showed

528 intermediate susceptibility to this aminoglycoside.

529       Presence of a gene in a bacterium's genome does not mean it is functionally active, nor does it

530 give any indication as to how active the gene is if it is indeed functional: e.g. all nine *K. pneumoniae*

isolates encoded KpnG and KpnH (strict CARD matches**; Figure 2a**), but only two showed any

resistance to gentamicin upon susceptibility testing.

Isolates #64 (SHV-1), #74 (SHV-36), #85, #88, #91 (OKP-A-2), #92 (SHV-36), #95 (SHV-28), #102 (SHV-164), #115 (SHV-36) and #118 (SHV-36) – which all encoded extended-spectrum β-lactamases (SHV) or *Escherichia ampH* β-lactamase (OKP-A-2) but lacked OmpK35 and OmpK36 – showed intermediate susceptibility to the carbapenem meropenem. Loss of the two porins OmpK35 and OmpK36 is known to confer resistance to carbapenems in strains producing extended-spectrum β-lactamases or plasmid-mediated AmpC-type β-lactamases (93). *K. pneumoniae* #64 was isolated from an infant with clinically diagnosed NEC with confirmed *Klebsiella* colonization. Importantly this preterm infant had previously been treated with benzylpenicillin, gentamicin and meropenem, which may link to the observed phenotypic resistance and corresponding AMR genes *Escherichia ampH* β-lactamase, $bla_{SHV-1}$ and KpnG/KpnH and SHV-1, and lack of OmpK35 and OmpK36, and suggests further treatment with gentamicin and meropenem would have been ineffective in this infant. Indeed, the infant was treated with cefotaxime, metronidazole and vancomycin in a subsequent round of medication (**Supplementary Table 2**). *K. pneumoniae* #115 was isolated from a baby that had confirmed NEC: the strain was resistant to benzylpenicillin (encoded *Escherichia ampH* β-lactamase and $bla_{SHV-36}$) and showed intermediate resistance to meropenem (lacked OmpK35 and OmpK36). *K. michiganensis* #88, also isolated from a baby that had NEC, showed intermediate resistance to both benzylpenicillin ($bla_{OXY-1-2}$, perfect CARD match; **Figure 3b**) and meropenem (lacked OmpK35 and OmpK36): both antibiotics had been administered to the baby at birth. *K. grimontii* #99, isolated from a baby with suspected sepsis, showed intermediate resistance to benzylpenicillin (encoded *Escherichia ampH* β-lactamase and $bla_{OXY-6-2}$). These data indicate that preterm-associated *Klebsiella* have a multi-drug-resistant phenotype that may prove problematic when treatment options are required for sepsis or NEC. Interestingly, other isolates (e.g. #95, recovered from an infant who had received benzylpenicillin and gentamicin; **Supplementary Table 2**) associated with 'healthy' preterm infants also harboured AMR genes (#95: **Figure 2a**) and phenotypic resistance profiles suggesting that administration of antibiotics to preterm infants with no signs of clinical infection contributes to

558    the reservoir of AMR genes – the 'resistome' (94) – which may increase horizontal gene transfer of

559    AMR determinants to other opportunistic pathogens residing within the GI tract.

560

561    **Abundance of *K. oxytoca* and related species in metagenomic datasets**

562    We used a published metagenomics dataset (1) to determine the prevalence of *K. oxytoca*, *K.*

563    *michiganensis* and *K. grimontii* in the preterm infant gut microbiome. These data had previously been

564    used to look at the relationship between NEC and uropathogenic *Esc. coli*, and metadata were

565    available for the samples. Ward *et al.* (1) collected a total of 327 samples at three stages of infant life:

566    stage1, days 3–9; stage2, days 10–16; stage3, days 17–22. Within each life stage samples were

567    collected on more than one day for some infants. In the current study, only samples processed under

568    Protocol A of Ward *et al.* (1) and from the earliest collection day within each life stage were analysed.

569    For those samples for which multiple sets of paired-end data were available, read data were

570    concatenated and used in analyses (**Supplementary Table 3**).

571    Stage1 comprised samples from 127 infants (105 preterm, 22 term), 16 of whom had been

572    diagnosed with NEC and 10 infants had subsequently died. Stage2 comprised samples from 146

573    infants (128 preterm and 18 term), 24 of whom later developed NEC with 18 deaths. Stage3

574    comprised samples from 54 infants (48 preterm, 6 term), including eight NEC patients, six of whom

575    died. Samples were collected from 165 distinct infants (143 preterm, 22 term) but only 41 of them

576    were sequenced at all three life stages. Infants were born either vaginally ($n = 70$) or by caesarean

577    section ($n = 95$). The gestational ages of preterm infants ranged from 23 to 29 weeks (mean 26.1

578    weeks), while the term babies ranged from 38 to 41 weeks (mean 39.2 weeks).

579    As we had found that the MetaPhlAn2.6 database contained non-*K. oxytoca* genomes within its

580    *K. oxytoca* dataset (detecting *K. oxytoca*, *K. michiganensis* and *R. ornithinolytica* (refer to Methods)),

581    we used Centrifuge to determine abundance of this species in metagenomes (**Figure 5ab**). Due to

582    their genomic similarity, *K. oxytoca* and *K. michiganensis* could not be readily distinguished using

583    Centrifuge (**Figure 5a**); no genomes assigned to *K. grimontii* were included in the Centrifuge

584    database at the time this study was undertaken. Though it should be noted that, while Centrifuge (and

585    Kraken2) relies on NCBI taxonomy for species identification, there are still many genomes within

586  GenBank/RefSeq that are assigned to the wrong species (e.g. assemblies GCA_001052235.1 and

587  GCA_000427015.1 within our curated pangenome dataset have been confirmed by detailed analyses

588  to be *K. grimontii* (**Supplementary Figure 2** and (36)), but still assigned as *K. oxytoca* and *K.*

589  *michiganensis*, respectively, within GenBank as of 28 July 2019; these are by no means the only

590  examples from our current study).

591     For those samples harbouring *K. oxytoca*, relative abundance of the bacterium increased from

592  stage1 to stage2 and decreased at stage3 (**Figure 5b**).

593

594  **Metapangenome analysis of preterm infant metagenomic data to detect *K. oxytoca*, *K.***

595  ***michiganensis* and *K. grimontii***

596     Using a set of 162 *K. oxytoca*-related genomes (**Supplementary Table 4**) and those of the five

597  infant isolates, a pangenome was generated using Roary. The pangenome dataset consisted of 76 *K.*

598  *michiganensis* (mean ANI among strains 98.55 (0.60) %, range 97.13–100 %), 64 *K. oxytoca* (mean

599  ANI among strains 99.20 (0.30) %, range 98.53–100 %) and 27 *K. grimontii* (mean ANI among

600  strains 98.45 (1.47) %, range 95.70–100 %) strains. A total of 40,605 genes were detected in the open

601  pangenome: 2,769 of them constituted the core gene cluster, while the accessory cluster included

602  5,108 genes and the remaining 32,728 genes formed the strain-specific cluster. A PCA plot based on

603  the accessory genes clustered strains into the three different species (**Figure 5c**), in agreement with

604  our phylogenetic analysis of the core genes (**Supplementary Figure 5b**) and consistent with the

605  findings of (76), who were able to split the three species (phylogroups) based on a pangenome

606  analysis of fewer genomes.

607     The Roary-generated pangenome was used as a custom database for PanPhlAn, to detect the

608  presence and absence of core genes and accessory genes in each infant sample. As expected, the

609  proportion of reads that PanPhlAn mapped to the custom database correlated with the Centrifuge-

610  generated abundance data (**Figure 5d**). In stage1, 13 infants (12 preterm, 1 term; 10.2 % of all stage1

611  samples) were predicted to carry *K. oxytoca*-related species (**Figure 5e**); in stage2, the number was 24

612  (22 preterm, 2 term; 16.4 % of all stage2 samples) (**Figure 5f**); in stage3 infants, the rate of carriage

613  was much lower, with only three infants (all preterm; 5.6 % of all stage3 samples) potentially

614    harbouring target species (**Figure 5g**). The change in prevalence of *K. oxytoca*-related species across

615    the three stages based on the PanPhlAn analysis was consistent with the *K. oxytoca* abundance data

616    generated with Centrifuge (**Figure 5b**).

617        The pangenome accessory genes were used in PCA to define which species strains detected by

618    PanPhlAn belonged to. In stage1, samples from six preterm infants (P10111, P10141, P10301,

619    P10451, P11292, P12121) and one term infant (P30221) harboured *K. michiganensis* (**Figure 5e**). In

620    stage2, samples from four preterm infants (P10471, P10472, P11311, P11701) harboured *K. oxytoca*,

621    while those from 14 preterm infants (P10071, P10231, P10301, P10441, P10451, P10501, P10601,

622    P11151, P11202, P11291, P11292, P12121, P12641, P12651) and one term infant (P30221)

623    harboured *K. michiganensis* (**Figure 5f**). Four samples (P11351, P12621, P20241, P30141) could not

624    be assigned a species, while the sample from P11981 located close to *K. oxytoca* (**Figure 5f**).

625    Similarly, two of the stage3 samples from P12221 and P12651 carried *K. michiganensis*, while

626    P11981 located near *K. oxytoca* (**Figure 5g**).

627

628    **Recovery of *K. oxytoca* and *K. michiganensis* MAGs from metagenomes**

629        Since the abundance of *K. oxytoca*-related species was considerable in some infant samples, we

630    attempted to obtain high-quality MAGs directly from these metagenomes and to assign them to the

631    species. The metagenomic samples were checked and their reads aligned against those of the 167-

632    genome database; reads that mapped were extracted and assembled as 'original MAGs'. The genome

633    sizes of the 167 genomes ranged from 5.72 Mb to 7.23 Mb (mean 6.35 Mb), thus a genome size of at

634    least 5.5 Mb was used to define a likely complete MAG. After assembly of the reads that mapped to

635    our database, MAGs were generated from the stage 1 (_s1), stage2 (_s2) and stage3 (_s3) samples.

636    ANI and phylogenetic analyses showed these MAGs to be *K. michiganensis* (P10301_s1, P10451_s1,

637    P11292_s1, P12121_s1, P30221_s1, P10071_s2, P10301_s2, P10441_s2, P10451_s2, P10501_s2,

638    P10601_s2, P11151_s2, P11202_s2, P11291_s2, P11292_s2, P12121_s2, P12641_s2, P12651_s2,

639    P30221_s2, P12221_s3 and P12651_s3; mean ANI with GCA_002925905 of 98.61 ± 0.66 %) or *K.*

640    *oxytoca* (P10472_s2, P11311_s2, P11701_s2, P11981_s2; mean ANI with GCA_900083895 of 99.30

641  ± 0.26 %) (36). No *K. grimontii* MAGs were recovered from any samples.

642  Prior to checking the completeness and contamination of the *K. michiganensis* and *K. oxytoca*

643  MAGs, contigs <500 nt in length were removed from the assemblies. A high-quality MAG requires a

644  >90 % genome completeness with contamination <5 % (44). According to CheckM results, three

645  stage1 MAGs (P10301_s1, P10451_s1, P30221_s1), seven stage2 MAGs (P10441_s2, P10451_s2,

646  P10501_s2, P10601_s2, P11291_s2, P11292_s2, P30221_s2) and one stage3 MAG (P12221_s3) were

647  of high quality. The rest of the MAGs were ≥90 % complete, but were contaminated (e.g. P12621_2

648  contained 274.59 % contamination). Thus, we attempted to decontaminate the MAGs using a

649  Diamond BLAST-based approach.

650  Since we already knew the species each MAG belonged to from PCA and PanPhlAn analyses,

651  scaffolds were mapped against the relevant species-specific genome database under different

652  minimum BLAST identity to report alignments, which were used to generate 'cleaner' MAGs.

653  **Supplementary Figure 6** shows the change in genome completeness and the percentage

654  contamination of stage2 MAGs when different blast identities were applied. The changes were

655  negligible for those MAGs with high-quality-level completeness and lacking contaminants even when

656  the cut-off was set at 99 %. For contaminated MAGs, the percentage contamination decreased

657  markedly as the BLAST identity became stricter and reduced to the bottom when all scaffolds in that

658  MAG could be aligned with 100 % identity. However, 100 % was not an appropriate threshold as the

659  genome completeness was affected greatly at this point (**Supplementary Figure 6**). Instead, a cut-off

660  of 99 % was used to decontaminate MAGs because 13 high-quality level MAGs and 5 medium-

661  quality level MAGs could be obtained when using this identity threshold. Stage2 MAGs that passed

662  PanPhlAn, PCA and ANI analysis reached at least reach medium-quality level using a 99 % identity

663  threshold. This cut-off was also suitable for stage1 MAGs, the quality of which was high. However,

664  for stage3 MAGs, the percentage contamination from P11981_s3 decreased to medium-quality level

665  only at 100 % identity, at which time the genome completeness fell down to 42.40 %. After

666  evaluating their genome completeness and contamination levels, a total of 25 MAGs (**Table 2**) were

667  assessed further.

668  The presence of tRNAs for the standard 20 amino acids and rRNA was examined as a secondary

669    measure of genome quality. A high-quality MAG requires at least 18 of the 20 possible amino acids

670    (44). P11981_s2 (16 aa) and P12651_s3 (17 aa) had to be classified as medium-quality MAGs even

671    though their genome completeness and contamination reached high-quality levels. 16S rRNA genes

672    were detected in all MAGs except P11151_s2, which was subsequently classified as medium quality.

673    Taking mandatory genome information into consideration (44), a total of 19 high-quality and six

674    medium-quality MAGs were recovered (**Table 2**); the sequences of these MAGs are available from

675    figshare. All of the MAGs had ≥15 standard tRNAs. High-quality MAGs had tRNAs that encoded an

676    average of 19.6 (0.7) of the 20 amino acids, some of them even had a tRNA that encodes an additional

677    amino acid SeC, while medium-quality MAGs had 18 (1.4) basic amino acids encoded by tRNAs.

678    High-quality MAGs consisted of ≤500 scaffolds in 52.6 % of cases (mean 600) and had an average

679    N50 of 121 kb, while only one medium-quality MAG comprised ≤500 scaffolds (mean 1174) and the

680    average N50 was less than half of that of high-quality MAGs (48.3 kb).

681

682    **Genotyping of MAGs**

683    Comparison of the sequences of the MAGs showed each infant harboured a different strain of *K.*

684    *oxytoca* (**Figure 6a**) or *K. michiganensis* (**Figure 6b**). In infants where MAGs were recovered across

685    different life stages, the MAGs were highly similar to one another (**Figure 6b**). Similar to what we

686    had seen with our isolates, the MAGs encoded a range of β-lactamase and virulence genes

687    (**Supplementary Figure 7**). It was also notable that two of the MAGs (*K. michiganensis* 10071_s2,

688    *K. oxytoca* 10472_s2) encoded *mcr-9* (perfect match), a plasmid-mediated colistin resistance gene and

689    phosphoethanolamine transferase. However, as noted above for our isolate work, presence of the

690    aforementioned genes in MAGs does not mean they were functionally active in the infants' GI tracts.

691    All the *K. michiganensis* MAGs encoded the siderophore enterobactin, along with all but one

692    (11981_s2) of the *K. oxytoca* MAGs. The allantoinase gene cluster associated with liver infection was

693    detected in the four *K. oxytoca* MAGs, but only a third of the *K. michiganensis* MAGs. We only

694    detected this cluster in the *K. grimontii* isolates we recovered (**Figure 3b**).

695    MLST analysis assigned 10 MAGs to eight known STs (**Table 2**). We believe insufficient

696    sequence coverage meant we were unable to ST more MAGs: e.g. *gapA* of MAG P11981_s2 aligned

697     exactly with *gapA* allele 2 sequence but it was only partial, leaving an incomplete match.

698         The STs of all genomes in the curated genome dataset were also identified (**Supplementary**

699     **Table 4**). Due to our limited understanding of *K. oxytoca*-related species, many combinations of

700     alleles have not been assigned corresponding STs yet, especially for the newly described species *K.*

701     *grimontii* (73). MLST identification showed that 54/64 *K. oxytoca sensu stricto* strains had known

702     STs, with some STs being more dominant than others. ST2, the most prevalent ST and represented by

703     15 strains, belonged to CC2. ST18 and ST19 were also in CC2, with both represented by two strains.

704     ST199 and ST176 were the second and third most common STs, respectively. Among *K.*

705     *michiganensis* strains, 48/76 of them could be assigned a known ST and comprised 21 distinct STs,

706     with nine represented by more than one isolate. ST11, ST27, ST50, ST85, ST143 and ST202 were the

707     most frequent STs, all of which were represented by at least four strains. *K. michiganensis* #108 was

708     ST157. Only 5/27 *K. grimontii* strains could be assigned an ST (#83 – ST72; #99 – ST76; 10-5250 –

709     ST47, GCA_000247915.1; 1148_KOXY – ST186, GCA_001052235.1; M5al – ST104,

710     GCA_001633115.1).

711

712     **SUMMARY**

713         *Klebsiella* spp. encode numerous virulence and antibiotic resistance genes that may contribute to

714     the pathogenesis of NEC and LOS. In this study, we characterized nine *K. pneumoniae*, three *K.*

715     *grimontii* and two *K. michiganensis* strains isolated prospectively from the faeces of a UK cohort of

716     preterm infants, and have shown these gut isolates are able to reside and persist in macrophages,

717     suggesting they can evade the immune system. These isolates will be used in future studies aiming to

718     replicate aspects of NEC and sepsis in model systems to confirm the role of *Klebsiella* spp. in these

719     diseases.

720         We have shown that mis-annotated genomes are being used in bioinformatics tools routinely

721     used to characterize the human gut microbiome. By using a carefully curated dataset to undertake

722     metapangenome analyses of the closely related species *K. oxytoca*, *K. michiganensis* and *K. grimontii*,

723     we have demonstrated that *K. michiganensis* is likely to be more clinically relevant to a subset of

724     preterm infants than *K. oxytoca*. Identity of publicly available genomes should be confirmed upon

725 download and linked to accurate taxonomic frameworks prior to analyses of data, especially when

726 attempting to identify and type closely related species in metagenomic data.

727

728 **AUTHOR STATEMENTS**

729 **Authors and contributions**

730     LH and LJH conceived and designed the study. All authors contributed to the writing of the

731 manuscript. TCB did the isolation and initial characterization work, and isolated DNA from bacteria.

732 YC and LH did bioinformatics associated with genome analyses. OL did the anvi'o work. TCB, YC

733 and OL were supervised by LH. LJH is chief investigator for the preterm clinical study from which

734 samples were used, and PC is the clinical lead for the study. CAG performed 16S rRNA gene library

735 preparation from faecal samples, and determined MIC values for strains. IO did macrophage assays.

736 CZS did iron assays, and additional siderophore bioinformatics analysis. SC performed the 16S rRNA

737 gene-associated bioinformatics, SP did 16S rRNA gene sequence analyses, and clinical database

738 management. LJH supervised CA, IO, CZS, SC and SP.

739

740 **Conflicts of interest**

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

742

753    and the Institute Strategic Programme Gut Microbes and Health BB/R012490/1, and its constituent

754    project(s) BBS/E/F/000PR10353 and BBS/E/F/000PR10356 to LJH.

755

756    **Ethics approval**

757        The Ethics Committee of the Faculty of Medical and Health Sciences of the University of

758    East Anglia (Norwich, UK) approved this study. The protocol for faeces collection was laid out by the

759    Norwich Research Park (NRP) Biorepository (Norwich, UK) and was in accordance with the terms of

760    the Human Tissue Act 2004 (HTA), and approved with licence number 11208 by the Human Tissue

761    Authority. Infants admitted to the NICU of the NNUH were recruited by doctors or nurses with

762    informed and written consent obtained from parents.

763

764    **Consent for publication**

765        All authors approved submission of the manuscript for publication.

766

774

775    **REFERENCES**

776    1.  Ward DV, Scholz M, Zolfo M, Taft DH, Schibler KR, Tett A, et al. Metagenomic sequencing

777        with strain-level resolution implicates uropathogenic *E. coli* in necrotizing enterocolitis and

778        mortality in preterm infants. Cell Rep. 2016 Mar 29;14(12):2912–24.

779    2.  Korpela K, de Vos WM. Early life colonization of the human gut: microbes matter

780        everywhere. Curr Opin Microbiol. 2018 Aug 4;44:70–8.

781  3.  Newburg DS, Walker WA. Protection of the neonate by the innate immune system of

782  developing gut and of human milk. Pediatr Res. 2007 Jan;61(1):2–8.

783  4.  Wandro S, Osborne S, Enriquez C, Bixby C, Arrieta A, Whiteson K. The microbiome and

784  metabolome of preterm infant stool are personalized and not driven by health outcomes,

785  including necrotizing enterocolitis and late-onset sepsis. mSphere. 2018 Jun 27;3(3).

786  5.  Wang Y, Hoenig JD, Malin KJ, Qamar S, Petrof EO, Sun J, et al. 16S rRNA gene-based

787  analysis of fecal microbiota from preterm infants with and without necrotizing enterocolitis.

788  ISME J. 2009 Aug;3(8):944–54.

789  6.  DiBartolomeo ME, Claud EC. The developing microbiome of the preterm infant. Clin Ther.

790  2016;38(4):733–9.

791  7.  Magne F, Abély M, Boyer F, Morville P, Pochart P, Suau A. Low species diversity and high

792  interindividual variability in faeces of preterm infants as revealed by sequences of 16S rRNA

793  genes and PCR-temporal temperature gradient gel electrophoresis profiles. FEMS Microbiol

794  Ecol. 2006 Jul;57(1):128–38.

795  8.  Mai V, Torrazza RM, Ukhanova M, Wang X, Sun Y, Li N, et al. Distortions in development

796  of intestinal microbiota associated with late onset sepsis in preterm infants. PloS One.

797  2013;8(1):e52876.

798  9.  Shaw AG, Sim K, Randell P, Cox MJ, McClure ZE, Li M-S, et al. Late-onset bloodstream

799  infection and perturbed maturation of the gastrointestinal microbiota in premature infants. PloS

800  One. 2015;10(7):e0132923.

801  10.  Warner BB, Tarr PI. Necrotizing enterocolitis and preterm infant gut bacteria. Semin Fetal

802  Neonatal Med. 2016;21(6):394–9.

803  11.  Robertson C, Savva G, Clapuci R, Jones J, Maimouni H, Brown E, et al. Incidence of

804  necrotising enterocolitis before and after introducing routine prophylactic *Lactobacillus and*

805  *Bifidobacterium* probiotics. Arch Child Fetal Neonatal Ed [Internet]. 2019 [cited 2019 Dec 3];

806  Available from: https://www.ncbi.nlm.nih.gov/pubmed/31666311

807  12.  Lu L, Claud EC. Intrauterine inflammation, epigenetics, and microbiome influences on

808  preterm infant health. Curr Pathobiol Rep. 2018;6(1):15–21.

809    13. Adeolu M, Alnajar S, Naushad S, S Gupta R. Genome-based phylogeny and taxonomy of the

810        'Enterobacteriales': proposal for Enterobacterales ord. nov. divided into the families

811        Enterobacteriaceae, Erwiniaceae fam. nov., Pectobacteriaceae fam. nov., Yersiniaceae fam.

812        nov., Hafniaceae fam. nov., Morganellaceae fam. nov., and Budviciaceae fam. nov. Int J Syst

813        Evol Microbiol. 2016 Dec;66(12):5575–99.

814    14. Pammi M, Cope J, Tarr PI, Warner BB, Morrow AL, Mai V, et al. Intestinal dysbiosis in

815        preterm infants preceding necrotizing enterocolitis: a systematic review and meta-analysis.

816        Microbiome. 2017 09;5(1):31.

817    15. Olm MR, Bhattacharya N, Crits-Christoph A, Firek BA, Baker R, Song YS, et al. Necrotizing

818        enterocolitis is preceded by increased gut bacterial replication, Klebsiella, and fimbriae-

819        encoding bacteria that may stimulate TLR4 receptors. bioRxiv. 2019 Feb 22;558676.

820    16. Claud EC, Keegan KP, Brulc JM, Lu L, Bartels D, Glass E, et al. Bacterial community

821        structure and functional contributions to emergence of health or necrotizing enterocolitis in

822        preterm infants. Microbiome. 2013 Jul 10;1(1):20.

823    17. Sim K, Shaw AG, Randell P, Cox MJ, McClure ZE, Li M-S, et al. Dysbiosis anticipating

824        necrotizing enterocolitis in very premature infants. Clin Infect Dis Off Publ Infect Dis Soc Am.

825        2015 Feb 1;60(3):389–97.

826    18. Torrazza RM, Ukhanova M, Wang X, Sharma R, Hudak ML, Neu J, et al. Intestinal microbial

827        ecology and environmental factors affecting necrotizing enterocolitis. PloS One.

828        2013;8(12):e83304.

829    19. Alcon-Giner C, Dalby MJ, Caim S, Ketskemety J, Shaw A, Sim K, et al. Microbiota

830        supplementation with Bifidobacterium and Lactobacillus modifies the preterm infant gut

831        microbiota and metabolome. bioRxiv. 2019 Jul 12;698092.

832    20. Kiu R, Caim S, Alcon-Giner C, Belteki G, Clarke P, Pickard D, et al. Preterm infant-

833        associated Clostridium tertium, Clostridium cadaveris, and Clostridium paraputrificum strains:

834        genomic and evolutionary insights. Genome Biol Evol. 2017 01;9(10):2707–14.

835    21. Harris SR, Feil EJ, Holden MTG, Quail MA, Nickerson EK, Chantratita N, et al. Evolution of

836        MRSA during hospital transmission and intercontinental spread. Science. 2010 Jan

837    22;327(5964):469–74.

838    22. Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson O, Huttenhower C. Metagenomic

839        microbial community profiling using unique clade-specific marker genes. Nat Methods. 2012

840        Aug;9(8):811–4.

841    23. Abbott JC. BugBuilder - an automated microbial genome assembly and analysis pipeline.

842        bioRxiv. 2017 Jun 11;148783.

843    24. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI

844        analysis of 90K prokaryotic genomes reveals clear species boundaries. Nat Commun. 2018

845        30;9(1):5114.

846    25. Jolley KA, Bray JE, Maiden MCJ. Open-access bacterial population genomics: BIGSdb

847        software, the PubMLST.org website and their applications. Wellcome Open Res. 2018;3:124.

848    26. Lam M, Wyres KL, Duchêne S, Wick RR, Judd LM, Gan Y-H, et al. Population genomics of

849        hypervirulent *Klebsiella pneumoniae* clonal-group 23 reveals early emergence and rapid global

850        dissemination. Nat Commun [Internet]. 2018 Jul 13 [cited 2018 Aug 14];9. Available from:

851        https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6045662/

852    27. Wick RR, Heinz E, Holt KE, Wyres KL. Kaptive Web: User-Friendly Capsule and

853        lipopolysaccharide serotype prediction for *Klebsiella* genomes. J Clin Microbiol. 2018

854        Jun;56(6).

855    28. Wyres KL, Wick RR, Gorrie C, Jenney A, Follador R, Thomson NR, et al. Identification of

856        *Klebsiella* capsule synthesis loci from whole genome data. Microb Genomics. 2016

857        Dec;2(12):e000102.

858    29. Chen L, Zheng D, Liu B, Yang J, Jin Q. VFDB 2016: hierarchical and refined dataset for big

859        data analysis--10 years on. Nucleic Acids Res. 2016 Jan 4;44(D1):D694-697.

860    30. Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P, Tsang KK, et al. CARD 2017:

861        expansion and model-centric curation of the comprehensive antibiotic resistance database.

862        Nucleic Acids Res. 2017 Jan 4;45(D1):D566–73.

863    31. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an

864        advanced analysis and visualization platform for 'omics data. PeerJ. 2015;3:e1319.

865     32. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic

866         gene recognition and translation initiation site identification. BMC Bioinformatics. 2010 Mar

867         8;11:119.

868     33. van Dongen S, Abreu-Goodger C. Using MCL to extract clusters from networks. Methods

869         Mol Biol Clifton NJ. 2012;804:281–95.

870     34. Schwyn B, Neilands JB. Universal chemical assay for the detection and determination of

871         siderophores. Anal Biochem. 1987 Jan;160(1):47–56.

872     35. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of

873         metagenomic sequences. Genome Res. 2016;26(12):1721–9.

874     36. Chen Y. Genome analysis of Gram-negative bacteria isolated from preterm baby faeces and

875         whole-genome analysis of *Klebsiella oxytoca* [MSc]. Imperial College London; 2018.

876     37. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinforma Oxf Engl. 2014 Jul

877         15;30(14):2068–9.

878     38. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary: rapid

879         large-scale prokaryote pan genome analysis. Bioinforma Oxf Engl. 2015 Nov 15;31(22):3691–

880         3.

881     39. Price MN, Dehal PS, Arkin AP. FastTree: computing large minimum evolution trees with

882         profiles instead of a distance matrix. Mol Biol Evol. 2009 Jul;26(7):1641–50.

883     40. Scholz M, Ward DV, Pasolli E, Tolio T, Zolfo M, Asnicar F, et al. Strain-level microbial

884         epidemiology and population genomics from shotgun metagenomics. Nat Methods.

885         2016;13(5):435–8.

886     41. Lê S, Josse J, Husson F. FactoMineR: an R package for multivariate analysis. J Stat Softw

887         [Internet]. 2008 [cited 2018 May 1];25(1). Available from:

888         https://www.jstatsoft.org/article/view/v025i01

889     42. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a

890         new genome assembly algorithm and its applications to single-cell sequencing. J Comput Biol J

891         Comput Mol Cell Biol. 2012 May;19(5):455–77.

892     43. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the

893 quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome

894 Res. 2015 Jul;25(7):1043–55.

895 44. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.

896 Minimum information about a single amplified genome (MISAG) and a metagenome-assembled

897 genome (MIMAG) of bacteria and archaea. Nat Biotechnol. 2017 Aug 8;35(8):725–31.

898 45. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. Nat

899 Methods. 2015 Jan;12(1):59–60.

900 46. Chen Y, Brook T, Alcon-Giner C, Clarke P, Hall L, Hoyles L. Draft genome sequences of

901 *Citrobacter freundii* and *Citrobacter murliniae* strains isolated from the feces of preterm infants.

902 Microbiol Resour Announc. 2019;8:e00494-19.

903 47. Chen Y, Brook T, Alcon-Giner C, Clarke P, Hall L, Hoyles L. Draft genome sequence of

904 *Raoultella ornithinolytica* P079F W, isolates from the feces of a preterm infant. Microbiol

905 Resour Announc. 2019;8:e00493-19.

906 48. Dobbler PT, Procianoy RS, Mai V, Silveira RC, Corso AL, Rojas BS, et al. Low microbial

907 diversity and abnormal microbial succession is associated with necrotizing enterocolitis in

908 preterm infants. Front Microbiol. 2017;8:2243.

909 49. Greenwood C, Morrow AL, Lagomarcino AJ, Altaye M, Taft DH, Yu Z, et al. Early empiric

910 antibiotic use in preterm infants is associated with lower bacterial diversity and higher relative

911 abundance of *Enterobacter*. J Pediatr. 2014 Jul;165(1):23–9.

912 50. Itani T, Ayoub Moubareck C, Melki I, Rousseau C, Mangin I, Butel M-J, et al. Preterm

913 infants with necrotising enterocolitis demonstrate an unbalanced gut microbiota. Acta Paediatr

914 Oslo Nor 1992. 2018 Jan;107(1):40–7.

915 51. Vergnano S, Menson E, Kennea N, Embleton N, Russell AB, Watts T, et al. Neonatal

916 infections in England: the NeonIN surveillance network. Arch Dis Child Fetal Neonatal Ed.

917 2011 Jan;96(1):F9–14.

918 52. Rodrigues C, Passet V, Rakotondrasoa A, Diallo TA, Criscuolo A, Brisse S. Description of

919 *Klebsiella africanensis* sp. nov., *Klebsiella variicola* subsp. *tropicalensis* subsp. nov. and

920 *Klebsiella variicola* subsp. *variicola* subsp. nov. Res Microbiol. 2019 May;170(3):165–70.

921    53. Holt KE, Wertheim H, Zadoks RN, Baker S, Whitehouse CA, Dance D, et al. Genomic

922         analysis of diversity, population structure, virulence, and antimicrobial resistance in *Klebsiella*

923         *pneumoniae*, an urgent threat to public health. Proc Natl Acad Sci U S A. 2015 Jul

924         7;112(27):E3574-3581.

925    54. Turton JF, Payne Z, Coward A, Hopkins KL, Turton JA, Doumith M, et al. Virulence genes

926         in isolates of *Klebsiella pneumoniae* from the UK during 2016, including among carbapenemase

927         gene-positive hypervirulent K1-ST23 and 'non-hypervirulent' types ST147, ST15 and ST383. J

928         Med Microbiol. 2018 Jan;67(1):118–28.

929    55. Turton JF, Perry C, Elgohari S, Hampton CV. PCR characterization and typing of *Klebsiella*

930         *pneumoniae* using capsular type-specific, variable number tandem repeat and virulence gene

931         targets. J Med Microbiol. 2010 May;59(Pt 5):541–7.

932    56. Fung CP, Hu BS, Chang FY, Lee SC, Kuo BI, Ho M, et al. A 5-year study of the

933         seroepidemiology of *Klebsiella pneumoniae*: high prevalence of capsular serotype K1 in Taiwan

934         and implication for vaccine efficacy. J Infect Dis. 2000 Jun;181(6):2075–9.

935    57. Nishino K, Senda Y, Yamaguchi A. CRP regulator modulates multidrug resistance of

936         *Escherichia coli* by repressing the *mdtEF* multidrug efflux genes. J Antibiot (Tokyo). 2008

937         Mar;61(3):120–7.

938    58. Nagakubo S, Nishino K, Hirata T, Yamaguchi A. The putative response regulator BaeR

939         stimulates multidrug resistance of *Escherichia coli* via a novel multidrug exporter system,

940         MdtABC. J Bacteriol. 2002 Aug;184(15):4161–7.

941    59. Doménech-Sánchez A, Hernández-Allés S, Martínez-Martínez L, Benedí VJ, Albertí S.

942         Identification and characterization of a new porin gene of *Klebsiella pneumoniae*: its role in

943         beta-lactam antibiotic resistance. J Bacteriol. 1999 May;181(9):2726–32.

944    60. Gunn JS, Lim KB, Krueger J, Kim K, Guo L, Hackett M, et al. PmrA-PmrB-regulated genes

945         necessary for 4-aminoarabinose lipid A modification and polymyxin resistance. Mol Microbiol.

946         1998 Mar;27(6):1171–82.

947    61. Bohn C, Bouloc P. The *Escherichia coli cmlA* gene encodes the multidrug efflux pump

948         Cmr/MdfA and is responsible for isopropyl-beta-D-thiogalactopyranoside exclusion and

949  spectinomycin sensitivity. J Bacteriol. 1998 Nov;180(22):6072–5.

62.  Srinivasan VB, Rajamohan G. KpnEF, a new member of the *Klebsiella pneumoniae* cell
     envelope stress response regulon, is an SMR-type efflux pump involved in broad-spectrum
     antimicrobial resistance. Antimicrob Agents Chemother. 2013 Sep;57(9):4449–62.

63.  Srinivasan VB, Singh BB, Priyadarshi N, Chauhan NK, Rajamohan G. Role of novel
     multidrug efflux pump involved in drug resistance in *Klebsiella pneumoniae*. PloS One.
     2014;9(5):e96288.

64.  Esterly JS, Richardson CL, Eltoukhy NS, Qi C, Scheetz MH. Genetic Mechanisms of
     antimicrobial resistance of *Acinetobacter baumannii*. Ann Pharmacother. 2011 Feb;45(2):218–
     28.

65.  Tärnberg M, Nilsson LE, Monstein H-J. Molecular identification of (bla)SHV, (bla)LEN and
     (bla)OKP beta-lactamase genes in *Klebsiella pneumoniae* by bi-directional sequencing of
     universal SP6- and T7-sequence-tagged (bla)SHV-PCR amplicons. Mol Cell Probes. 2009
     Aug;23(3–4):195–200.

66.  Regué M, Climent N, Abitiu N, Coderch N, Merino S, Izquierdo L, et al. Genetic
     Characterization of the *Klebsiella pneumoniae waa* gene cluster, involved in core
     lipopolysaccharide biosynthesis. J Bacteriol. 2001 Jun;183(12):3564–73.

67.  Follador R, Heinz E, Wyres KL, Ellington MJ, Kowarik M, Holt KE, et al. The diversity of
     *Klebsiella pneumoniae* surface polysaccharides. Microb Genomics [Internet]. 2016 Aug 25
     [cited 2020 Mar 3];2(8). Available from:
     https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5320592/

68.  Beaugerie L, Metz M, Barbut F, Bellaiche G, Bouhnik Y, Raskine L, et al. *Klebsiella oxytoca*
     as an agent of antibiotic-associated hemorrhagic colitis. Clin Gastroenterol Hepatol Off Clin
     Pract J Am Gastroenterol Assoc. 2003 Sep;1(5):370–6.

69.  Högenauer C, Langner C, Beubler E, Lippe IT, Schicho R, Gorkiewicz G, et al. *Klebsiella
     oxytoca* as a causative organism of antibiotic-associated hemorrhagic colitis. N Engl J Med.
     2006 Dec 7;355(23):2418–26.

70.  Raveh-Sadka T, Thomas BC, Singh A, Firek B, Brooks B, Castelle CJ, et al. Gut bacteria are

977  rarely shared by co-hospitalized premature infants, regardless of necrotizing enterocolitis

978  development. eLife. 2015 Mar 3;4.

979  71.  Brisse S, Verhoef J. Phylogenetic diversity of *Klebsiella pneumoniae* and *Klebsiella oxytoca*

980  clinical isolates revealed by randomly amplified polymorphic DNA, gyrA and parC genes

981  sequencing and automated ribotyping. Int J Syst Evol Microbiol. 2001 May;51(Pt 3):915–24.

982  72.  Herzog KAT, Schneditz G, Leitner E, Feierl G, Hoffmann KM, Zollner-Schwetz I, et al.

983  Genotypes of *Klebsiella oxytoca* isolates from patients with nosocomial pneumonia are distinct

984  from those of isolates from patients with antibiotic-associated hemorrhagic colitis. J Clin

985  Microbiol. 2014 May;52(5):1607–16.

986  73.  Passet V, Brisse S. Description of *Klebsiella grimontii* sp. nov. Int J Syst Evol Microbiol.

987  2018 Jan;68(1):377–81.

988  74.  Saha R, Farrance CE, Verghese B, Hong S, Donofrio RS. *Klebsiella michiganensis* sp. nov., a

989  new bacterium isolated from a tooth brush holder. Curr Microbiol. 2013 Jan;66(1):72–8.

990  75.  Merla C, Rodrigues C, Passet V, Corbella M, Thorpe HA, Kallonen TVS, et al. Description of

991  *Klebsiella spallanzanii* sp. nov. and of *Klebsiella pasteurii* sp. nov. Front Microbiol.

992  2019;10:2360.

993  76.  Moradigaravand D, Martin V, Peacock SJ, Parkhill J. Population structure of multidrug-

994  resistant *Klebsiella oxytoca* within hospitals across the United Kingdom and Ireland identifies

995  sharing of virulence and resistance genes with *K. pneumoniae*. Genome Biol Evol. 2017 Mar

996  1;9(3):574–84.

997  77.  Hu Y, Wei L, Feng Y, Xie Y, Zong Z. *Klebsiella huaxiensis* sp. nov., recovered from human

998  urine. Int J Syst Evol Microbiol. 2019 Feb;69(2):333–6.

999  78.  Shibu P. Investigations of carbapenem-resistant *Klebsiella* species and associated clinical

1000  considerations. [London]: Westminster; 2019.

1001  79.  Chun J, Oren A, Ventosa A, Christensen H, Arahal DR, da Costa MS, et al. Proposed minimal

1002  standards for the use of genome data for the taxonomy of prokaryotes. Int J Syst Evol Microbiol.

1003  2018 Jan;68(1):461–6.

1004  80.  Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. DNA-

DNA hybridization values and their relationship to whole-genome sequence similarities. Int J Syst Evol Microbiol. 2007 Jan;57(Pt 1):81–91.

81. Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. Proc Natl Acad Sci U S A. 2009 Nov 10;106(45):19126–31.

82. Zheng B, Xu H, Yu X, Lv T, Jiang X, Cheng H, et al. Identification and genomic characterization of a KPC-2-, NDM-1- and NDM-5-producing *Klebsiella michiganensis* isolate. J Antimicrob Chemother. 2018 01;73(2):536–8.

83. Pedersen T, Sekyere JO, Govinden U, Moodley K, Sivertsen A, Samuelsen Ø, et al. Spread of plasmid-encoded NDM-1 and GES-5 Carbapenemases among extensively drug-resistant and pandrug-resistant clinical *Enterobacteriaceae* in Durban, South Africa. Antimicrob Agents Chemother. 2018;62(5).

84. Seiffert SN, Wüthrich D, Gerth Y, Egli A, Kohler P, Nolte O. First clinical case of KPC-3-producing *Klebsiella michiganensis* in Europe. New Microbes New Infect. 2019 May;29:100516.

85. Chou H-C, Lee C-Z, Ma L-C, Fang C-T, Chang S-C, Wang J-T. Isolation of a chromosomal region of *Klebsiella pneumoniae* associated with allantoin metabolism and liver infection. Infect Immun. 2004 Jul;72(7):3783–92.

86. Fevre C, Passet V, Weill F-X, Grimont PAD, Brisse S. Variants of the Klebsiella pneumoniae OKP chromosomal beta-lactamase are divided into two main groups, OKP-A and OKP-B. Antimicrob Agents Chemother. 2005 Dec;49(12):5149–52.

87. Cano V, March C, Insua JL, Aguiló N, Llobet E, Moranta D, et al. *Klebsiella pneumoniae* survives within macrophages by avoiding delivery to lysosomes. Cell Microbiol. 2015 Nov;17(11):1537–60.

88. Skaar EP. The battle for iron between bacterial pathogens and their vertebrate hosts. PLoS Pathog. 2010 Aug 12;6(8):e1000949.

89. Flo TH, Smith KD, Sato S, Rodriguez DJ, Holmes MA, Strong RK, et al. Lipocalin 2 mediates an innate immune response to bacterial infection by sequestering iron. Nature. 2004 Dec 16;432(7019):917–21.

1033    90.  Bachman MA, Oyler JE, Burns SH, Caza M, Lépine F, Dozois CM, et al. *Klebsiella*

1034         *pneumoniae* yersiniabactin promotes respiratory tract infection through evasion of lipocalin 2.

1035         Infect Immun. 2011 Aug;79(8):3309–16.

1036    91.  Lawlor MS, O'connor C, Miller VL. Yersiniabactin is a virulence factor for *Klebsiella*

1037         *pneumoniae* during pulmonary infection. Infect Immun. 2007 Mar;75(3):1463–72.

1038    92.  Navon-Venezia S, Kondratyeva K, Carattoli A. *Klebsiella pneumoniae*: a major worldwide

1039         source and shuttle for antibiotic resistance. FEMS Microbiol Rev. 2017 01;41(3):252–75.

1040    93.  Hamzaoui Z, Ocampo-Sosa A, Fernandez Martinez M, Landolsi S, Ferjani S, Maamar E, et

1041         al. Role of association of OmpK35 and OmpK36 alteration and blaESBL and/or blaAmpC genes

1042         in conferring carbapenem resistance among non-carbapenemase-producing *Klebsiella*

1043         *pneumoniae*. Int J Antimicrob Agents. 2018 Dec;52(6):898–905.

1044    94.  Penders J, Stobberingh EE, Savelkoul PHM, Wolffs PFG. The human microbiome as a

1045         reservoir of antimicrobial resistance. Front Microbiol. 2013;4:87.

1046

1047

**FIGURE LEGENDS**

1049 **Figure 1.** Summary information for UK cohort included in this study. (a) Breakdown of birth mode

1050 and sex of preterm neonates. (b) 16S rRNA gene sequence results for *Enterobacteriaceae*-positive

1051 samples: upper panel, samples from which lactose-positive isolates were recovered; lower panel,

1052 samples from which lactose-negative isolates were recovered. (c) Representation of lactose-negative

1053 and lactose-positive *Enterobacteriaceae* isolated from faecal samples. (d) Tentative identities of

1054 lactose-positive *Enterobacteriaceae* as determined by using API 20E.

1055

1056 **Figure 2.** Summary of (a) antibiotic resistance and (b) virulence factor genes found in the *K.*

1057 *pneumoniae* isolates by comparison of protein sequences with those of the CARD and VFDB,

1058 respectively. (a) Strict CARD match, not identical but the bit-score of the matched sequence is greater

1059 than the curated BLASTP bit-score cut-off; perfect CARD match, 100 % identical to the reference

1060 sequence along its entire length. Loose matches are not shown to avoid presenting false positives

1061 based on sequences with low homology and bit-scores below CARD BLASTP cut-off

1062 recommendations. (b) Identity (%), BLASTP reported only for those proteins sharing >70 % identity

1063 and 90 % query coverage with VFDB protein sequences.

1064

1065 **Figure 3.** Genomic characterization of the *K. michiganensis* and *K. grimontii* isolates recovered from

1066 neonates. (a) Anvi'o representation of the genomes of *K. michiganensis* and *K. grimontii* isolates

1067 recovered from different infants. It is clear the isolates are different from one another at the genomic

1068 level. (b) Virulence factor (left side) and antibiotic resistance (right side) genes encoded by the

1069 isolates. Criteria for identity and strict/perfect match with respect to VFDB and CARD, respectively,

1070 are the same as those given for **Figure 2**.

1071

1072 **Figure 4.** Phenotypic assays for the *Klebsiella* isolates recovered from infants. (a) Strains were tested

1073 for persistence in PMA-differentiated THP-1 macrophages using a gentamicin protection assay.

1074 Intracellular bacteria were enumerated 1.5 h and 6 h after infection to determine persistence (*n*=4).

1075 Results are shown as mean (SD). (b) *Klebseilla* strains were grown in minimal medium and at 20 h

1076     bacterial growth ($OD_{600}$) siderophore production was measured using the CAS assay ($n = 3$). Results

1077     are shown as mean (SD).

1078

1079     **Figure 5.** Identification of *K. oxytoca*-related species in infant faecal metagenomes. (a) Comparison

1080     of *K. oxytoca* and *K. michiganensis* abundance (as determined using Centrifuge) in stage2 samples of

1081     Ward *et al.* (1). (b) Abundance of *K. oxytoca* (determined using Centrifuge) across stage1, stage2 and

1082     stage3 samples of Ward *et al.* (1). (c) Separation of the strains of *K. grimontii* ($n$=27), *K.*

1083     *michiganensis* ($n$=76) and *K. oxytoca* ($n$=64) based on accessory genes ($n = 5,108$) detected in the

1084     Roary-generated open pangenome. (d) Relationship between PanPhlAn (overall alignment rate) and

1085     Centrifuge (abundance *K. oxytoca* (%)) data. (e, f, g) PCA plots show separation of strains in the

1086     pangenome plus PanPhlAn-detected strains based on presence of 500 randomly sampled accessory

1087     genes at (e) stage1, (f) stage2 and (g) stage3 of Ward *et al.* (1).

1088

1089     **Figure 6.** Anvi'o representation of the MAGs recovered from the metagenomes of infants included in

1090     the study of Ward *et al.* (1). (a) *K. oxytoca*. (b) *K. michiganensis*. It is notable that MAGs recovered

1091     from different life stages from the same infant (e.g. 10301_s1, 10301_s2) are highly similar to one

1092     another.

1093

1094  **Table 1.** Determination of MICs for *Klebsiella* spp. isolates

| Isolate ID | Species | Gentamicin (mg/L)* | Meropenem (mg/L)† | Benzylpenicillin (mg/L)‡ |
|---|---|---|---|---|
| #64 | *K. pneumoniae* | **6.25#** | <u>3.13</u> | 1560 |
| #74 | *K. pneumoniae* | 1.5625 | <u>6.25</u> | 3130 |
| #85 | *K. pneumoniae* | 1.5625 | <u>3.13</u> | 3130 |
| #91 | *K. quasipneumoniae* | 1.5625 | <u>6.25</u> | 3130 |
| #92 | *K. pneumoniae* | 1.5625 | <u>3.13</u> | 3130 |
| #95 | *K. pneumoniae* | <u>3.125</u> | <u>3.13</u> | 3130 |
| #102 | *K. pneumoniae* | 1.5625 | <u>6.25</u> | 3130 |
| #115 | *K. pneumoniae* | 1.5625 | <u>3.13</u> | 3130 |
| #118 | *K. pneumoniae* | 1.5625 | <u>6.25</u> | 3130 |
| #80 | *K. grimontii* | <u>3.125</u> | 1.56 | 6.25 |
| #83 | *K. grimontii* | **12.5** | 1.56 | 780 |
| #99 | *K. grimontii* | <u>3.125</u> | 0.78 | 780 |
| #88 | *K. michiganensis* | <u>3.125</u> | 3.13 | 3130 |
| #108 | *K. michiganensis* | **6.25** | 1.56 | 3130 |

1095  *\*Enterobacteriaceae* EUCAST breakpoint for gentamicin resistance is >4 mg/L, and for sensitivity is

1096  ≤2 mg/L.

1097  †*Enterobacteriaceae* EUCAST breakpoint for meropenem is resistance >8 mg/L, and for sensitivity is

1098  ≤2 mg/L.

1099  ‡No *Enterobacteriaceae* EUCAST data are available for benzylpenicillin.

1100  #Bold type, resistant; underlined, intermediate.

**Table 2.** Summary statistics for MAGs recovered from preterm infant metagenomes

| MAG* | Genome length (bp) | Max. contig length | Coverage† | N50 | No. of scaffolds | GC content (%) | Completeness (%)‡ | Contamination (%)‡ | CDS | No. of tRNAs | No. of rRNAs | Species | Quality | *gapA* | *infB* | *mdh* | *pgi* | *phoE* | *rpoB* | *tonB* | ST |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 10071_s2 | 6,577,866 | 334,901 | ~131x | 70,229 | 349 | 52.95 | 97.59 | 3.64 | 6204 | 54 | 3 | *K. michiganensis* | High | 3 | 9 | 8 | 9 | 20 | * | 8 | * |
| 10301_s1 | 6,304,211 | 308,355 | ~158x | 143,120 | 176 | 53.15 | 99.70 | 0.48 | 5920 | 54 | 4 | *K. michiganensis* | High | 3 | 5 | 21 | 13 | 74 | 6 | 12 | 202 |
| 10301_s2 | 6,504,340 | 160,223 | ~21x | 30,917 | 609 | 54.93 | 98.69 | 4.21 | 6191 | 50 | 2 | *K. michiganensis* | Medium | 3 | 5 | 21 | 13 | 74 | 6 | 12 | 202 |
| 10441_s2 | 6,128,395 | 396,555 | ~85x | 122,210 | 176 | 53.35 | 100.00 | 0.71 | 5691 | 58 | 4 | *K. michiganensis* | High | 3 | 5 | 21 | 13 | 24 | 6 | * | * |
| 10451_s1 | 6,135,593 | 499,885 | ~79x | 130,343 | 154 | 53.44 | 100.00 | 0.71 | 5708 | 53 | 7 | *K. michiganensis* | High | 3 | 5 | 21 | 13 | 24 | 6 | * | * |
| 10451_s2 | 6,140,759 | 377,549 | ~65x | 130,433 | 166 | 53.53 | 100.00 | 0.71 | 5707 | 54 | 4 | *K. michiganensis* | High | 3 | 5 | 21 | 13 | 24 | 6 | * | * |
| 10472_s2 | 6,352,656 | 271,779 | ~46x | 80,647 | 252 | 53.26 | 99.90 | 2.59 | 5889 | 50 | 3 | *K. oxytoca* | High | 1 | 7 | 2 | 1 | 65 | 1 | 2 | 176 |
| 10501_s2 | 6,179,579 | 451,918 | ~39x | 133,209 | 217 | 54.01 | 100.00 | 1.34 | 5730 | 55 | 9 | *K. michiganensis* | High | 3 | 5 | 21 | 3 | 24 | 6 | * | * |
| 10601_s2 | 6,014,627 | 328,826 | ~59x | 181,417 | 121 | 52.63 | 99.96 | 0.34 | 5540 | 56 | 4 | *K. michiganensis* | High | 3 | 5 | 21 | 3 | 24 | 6 | * | * |
| 11151_s2 | 6,595,368 | 114,676 | ~40x | 17,470 | 949 | 54.71 | 97.17 | 5.81 | 6252 | 59 | 0 | *K. michiganensis* | Medium | 3 | 5 | 21 | 13 | 20 | * | 12 | * |
| 11202_s2 | 6,328,179 | 267,673 | ~144x | 111,043 | 224 | 53.38 | 98.81 | 3.22 | 5864 | 54 | 3 | *K. michiganensis* | High | * | 8 | 24 | 33 | 20 | 6 | 23 | * |
| 11291_s2 | 6,301,261 | 429,108 | ~45x | 146,444 | 193 | 53.10 | 99.70 | 0.79 | 5911 | 68 | 8 | *K. michiganensis* | High | 3 | 8 | 17 | 21 | 40 | 17 | 29 | 84 |
| 11292_s1 | 6,533,348 | 267,672 | ~43x | 118,981 | 400 | 53.17 | 98.81 | 2.77 | 6064 | 48 | 3 | *K. michiganensis* | High | * | 8 | 24 | 33 | 20 | 6 | 23 | * |
| 11292_s2 | 6,383,200 | 267,500 | ~70x | 119,207 | 202 | 50.77 | 99.70 | 2.49 | 5891 | 63 | 8 | *K. michiganensis* | High | * | 8 | 24 | 33 | 20 | 6 | 23 | * |
| 11311_s2 | 6,427,314 | 266,393 | ~137x | 66,881 | 516 | 52.89 | 99.97 | 2.96 | 5986 | 46 | 1 | *K. oxytoca* | Medium | 2 | 2 | 2 | 3 | 19 | 2 | 2 | 199 |
| 11701_s2 | 6,230,855 | 162,114 | ~60x | 50,736 | 555 | 54.40 | 99.85 | 3.20 | 5807 | 50 | 2 | *K. oxytoca* | Medium | 1 | 7 | 2 | 1 | 65 | 1 | 2 | 176 |
| 11981_s2 | 5,658,923 | 237,765 | ~10x | 60,298 | 326 | 53.62 | 89.43 | 2.20 | 5280 | 30 | 2 | *K. oxytoca* | Medium | * | 2 | 2 | 3 | 19 | * | 2 | * |
| 12121_s1 | 5,843,313 | 426,086 | ~121x | 101,108 | 251 | 54.02 | 96.40 | 2.49 | 5358 | 53 | 4 | *K. michiganensis* | High | 3 | 33 | 17 | 45 | 20 | 6 | 48 | 149 |
| 12121_s2 | 5,690,611 | 180,068 | ~82x | 56,100 | 371 | 54.36 | 94.21 | 2.85 | 5226 | 35 | 3 | *K. michiganensis* | High | 3 | 33 | 17 | 45 | 20 | * | 48 | * |
| 12221_s3 | 6,115,156 | 469,104 | ~17x | 174,402 | 147 | 54.67 | 99.70 | 1.60 | 5667 | 63 | 8 | *K. michiganensis* | High | 3 | 5 | 21 | 20 | 24 | 6 | 30 | 108 |
| 12641_s2 | 6,383,005 | 237,677 | ~55x | 80,508 | 279 | 53.91 | 99.44 | 3.15 | 6027 | 56 | 2 | *K. michiganensis* | High | 3 | 5 | 21 | 13 | 74 | * | 12 | * |
| 12651_s2 | 5,124,388 | 199,319 | ~28x | 81,423 | 207 | 53.59 | 94.05 | 3.21 | 4801 | 42 | 3 | *K. michiganensis* | High | 14 | 24 | 15 | 8 | 18 | * | 4 | * |
| 12651_s3 | 5,422,817 | 363,696 | ~10x | 73,253 | 207 | 53.78 | 95.49 | 2.40 | 5047 | 41 | 2 | *K. michiganensis* | Medium | 14 | 2 | 15 | 8 | 18 | * | 4 | * |
| 30221_s1 | 6,220,970 | 480,926 | ~51x | 173,840 | 132 | 52.47 | 99.70 | 1.42 | 5770 | 61 | 9 | *K. michiganensis* | High | 3 | 5 | 21 | 20 | 11 | 6 | 20 | 43 |
| 30221_s2 | 6,224,584 | 349,605 | ~55x | 170,314 | 142 | 52.27 | 99.70 | 1.39 | 5765 | 58 | 7 | *K. michiganensis* | High | 3 | 5 | 21 | 20 | 11 | 6 | 20 | 43 |

1102    *E.g. 10301_s1 represents a MAG recovered from infant 10301 at stage1.

1103    †Coverage, based on coverage of longest contig (determined from spades data).

1104    ‡Completeness and contamination determined using CheckM (v1.0.18).