1    **Genome reconstruction and characterisation of extensively drug-resistant bacterial**

2    **pathogens through direct metagenomic sequencing of human faeces**

3

4    Andre Mu[1,2,*], Jason C. Kwong[1,2,3,*], Nicole S. Isles[1], Anders Gonçalves da Silva[1,2], Mark B.

5    Schultz[1,2], Susan A. Ballard[1,2], Glen P. Carter[2], Deborah A. Williamson[1,2], Torsten Seemann[2,4],

6    Timothy P. Stinear#[2], Benjamin P. Howden#[1,2,3]

7

8    [1]Microbiological Diagnostic Unit Public Health Laboratory, Department of Microbiology and

9    Immunology at the Peter Doherty Institute for Infection and Immunity, University of

10   Melbourne, Australia

11   [2]Doherty Applied Microbial Genomics, Department of Microbiology and Immunology at the

12   Peter Doherty Institute for Infection and Immunity, University of Melbourne, Australia

13   [3]Department of Infectious Diseases, Austin Health, Australia

14   [4]Melbourne Bioinformatics, University of Melbourne, Australia

15

16   *Authors contributed equally to the manuscript

17   # Authors contributed equally to the manuscript

18

19   **Correspondence:**

20   Professor Benjamin P. Howden

21   Director, Microbiological Diagnostic Unit Public Health Laboratory,

22   University of Melbourne at The Peter Doherty Institute for Infection and Immunity

23   Email: bhowden@unimelb.edu.au

24

25

26

27

28 **Abstract**

29 Whole-genome sequencing of microbial pathogens is revolutionising modern approaches to

30 outbreaks of infectious diseases and is reliant upon organism culture. Culture-independent

31 methods have shown promise in identifying pathogens, but high level reconstruction of

32 microbial genomes from microbiologically complex samples for more in-depth analyses

33 remains a challenge. Here, using metagenomic sequencing of a human faecal sample and

34 analysis by tetranucleotide frequency profiling projected onto emergent self-organising

35 maps, we were able to reconstruct the underlying populations of two extensively-drug

36 resistant pathogens, *Klebsiella pneumoniae* carbapenemase (KPC)-producing *Klebsiella*

37 *pneumoniae* and vancomycin-resistant *Enterococcus faecium.* From these genomes, we were

38 able to ascertain molecular typing results, such as MLST, and identify highly discriminatory

39 mutations in the metagenome to distinguish closely related strains. These proof-of-principle

40 results demonstrate the utility of clinical sample metagenomics to recover sequences of

41 important drug-resistant bacteria and application of the approach in outbreak investigations,

42 independent of the need to culture the organisms.

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

**Introduction**

Clinical and public health microbiology is undergoing a major transformation driven largely by high-throughput microbial genome sequencing. The application of microbial genomics in these areas has been well described, including use for high resolution microbial characterization, source and transmission tracking for nosocomial and community pathogens, and antimicrobial resistance detection and prediction [1-4]. Clinical and public health genomics, however, currently relies on routine culture-based assays to isolate pathogens of interest prior to whole genome sequencing, which presents inherent biases in analyses, and does not allow characterization of pathogens which are unculturable, below the level of culture detection, or unsuspected in a clinical sample [5-9].

Genomics-based approaches that overcome these difficulties by direct characterization of pathogens from clinical samples would be a major advance in clinical and public health microbiology. Metagenomics approaches can complement culture-based techniques for typing, and detecting, antimicrobial resistant (AMR) genes, and SNP variants [10,11]. Metagenomics allows for the sequencing of whole community genomic material extracted directly from clinical samples, such as faeces, blood, cerebrospinal fluid, sputum, and bronchoalveolar lavage fluid [12,13]. Current literature on public health metagenomics is largely based on interrogating the metagenome at 'first-order' level analyses; that is, either, characterizing bacterial biodiversity at the 16S rRNA gene level, identifying the functional profile of the microbial community, or alignment-based reference analyses for recovery of genomes [12]. However, the challenge still remains in applying alignment-free analyses of metagenomic data to obtain strain-level resolution that might help understand transmission of pathogens in a clinical setting. Recent metagenomics advances in the field of environmental microbiology and ecology may provide potential solutions here, including techniques that bin contigs based on their tetranucleotide frequency profiles (c.f.,[14-17]).

In this proof-of-principle study, we used whole community metagenomics and pathogen genome reconstruction to interrogate the metagenome of a patient colonised with an extensively drug-resistant pathogen KPC-producing *K. pneumoniae*. Here we demonstrate that faeces metagenomics not only identified detailed SNP information to distinguish clonal *K. pneumoniae* isolates, but also uncovered unsuspected colonization with vancomycin-

3

92    resistant *Enterococcus faecium* (VREfm), another high-risk antimicrobial resistant pathogen.

93    Furthermore, we investigated whether metagenomic analysis could be used to characterise

94    the resistance-harbouring genomes of a patient with long-term carriage of KPC-producing

95    Klebsiella pneumoniae, and link that patient to a local transmission network, independently

96    of the need for bacterial culture.

97

98

99 **Materials and Methods**

100 **Epidemiological context**

101 One faecal sample collected from a patient (Patient A) with known KPC-producing *K.*

102 *pneumoniae* colonisation underwent whole community metagenomics. Patient A was a

103 resident of an aged care facility and did not report any recent travel in the context of his

104 comorbidities and frailty, but had been an inpatient at a tertiary hospital with a known

105 outbreak of KPC-producing *K. pneumoniae* two months prior to sample collection.

106

107 **Whole community genomic DNA extraction and high throughput metagenomic sequencing**

108 Whole community gDNA was extracted from 0.2 grams of Patient faeces (herein referred to

109 as AUSMDU00008155) using the QiaAMP Stool kit following manufacturer's protocol with a

110 preprocessing step of mechanical lysis (Bertin Technologies precellys 24). MP Biomedicals'

111 Lysing Matrix B 2-ml tubes containing 0.1 mm silica beads were used for two 40 second cycles

112 of mechanical lysis (Bertin Technologies precellyis 24) at 6000x units with a 60 second rest on

113 ice in between. Genomic DNA from the faeces, and a no-template control, were processed

114 for sequencing using the Nextera XT kit on the Illumina MiSeq machine (V3, 600 cycles)

115 (Illumina Inc, San Diego, US) following a modified manufacturer's protocol. The following

116 modifications were included: a 1% (v/v) spike-in ratio of PhiX, denatured DNA was diluted to

117 a final concentration of 14.25 pM with pre-chilled HT1 buffer, and Tris-Cl 10 mM 0.1% Tween

118 20 was substituted with Qiagen's EB solution to dilute sequencing libraries and PhiX

119 throughout the protocol.

120

121 **Culture-dependent whole genome sequencing**

122 Concurrently, 14 individual colonies were picked at random from the same sample plated on

123 Brilliance™ CRE selective media (Thermo Fisher Scientific, Waltham, US), with each colony

124 undergoing whole-genome sequencing on the Illumina NextSeq 500 (Illumina Inc, San Diego,

125 US). Two KPC isolates (herein referred to as AUSMDU00008118 and AUSMDU00008119)

126 taken two days apart from a different patient that had previously undergone long-read

127 sequencing to investigate the plasmid dynamics within an outbreak of KPC-producing

128 *Klebsiella pneumoniae* were used as reference genomes for the analysis. DNA extraction, size

129 selection, and sequencing on the Pacific Biosciences RS II (Pacific Biosciences, Menlo Park, US)

130 were performed as previously described [18]. Genomic DNA from these isolates was also

5

131    sequenced on the Illumina NextSeq 500 for polishing to produce high quality closed genomes.

132

133    **Bioinformatic analyses**

134    **Metagenomic sequence data processing**

135    Metagenomic data from AUSMDU00008155 were processed prior to analysis with

136    Trimmomatic (v0.33) [19] for quality control and to remove adaptor sequences, PhiX

137    contamination, and trace contaminants from Illumina preparation kits. Paired-end reads were

138    merged and assembled using Iterative de Bruijn Graph De Novo Assembler for Uneven

139    sequencing Depth (IDBA-UD) [20] compiled for long reads (i.e., 651 bp). Further quality control

140    included removing host-derived gDNA using DeconSeq and the Human Genome Reference

141    Sequence (build 38; GCA_000001405.22) prior to downstream analyses.

142

143    **Metagenomic binning**

144    To reconstruct isolate genomes from the gut microbial community, an emergent self-

145    organizing map (ESOM) was used. Tetranucleotide frequencies were calculated for the

146    assembled contigs using Perl scripts developed by Dick et al., (2009) [16] in preparation for

147    analysis using ESOMs. The primary map structure was determined using in silico fragmented

148    (> 5kb) contigs; while contigs between 2.5kb and 5kb in length were projected onto the ESOM

149    using their tetranucleotide frequency profiles. Genomic binning was analysed using

150    Databionic ESOM Tool with default settings except K-Batch training algorithm in 200x400

151    windows, a starting value of 50 for the radius, and data points were normalized by RobustZT

152    transformation. Contigs with a native size smaller than 2.5kb were removed from analyses.

153    Reference genomes were included in the analysis to guide identification of "binned"

154    genomes, and validate completeness of genome recovery.

155

156    **Detection of antimicrobial resistance genes**

157    Assembled metagenomic contigs were screened for the presence of antimicrobial and

158    virulence genes, including carbapenemases ($bla_{KPC}$), using ABRicate

159    (https://github.com/tseemann/abricate). Briefly, ABRicate detects acquired resistance genes

160    using BLAST+ against the Resfinder database (Center for Genomic Epidemiology, University

161    of Denmark[21]).

162

6

163 **_In silico_ molecular typing and detection of antimicrobial resistance genes**

164 Assembled metagenomic contigs were screened for multi-locus sequence typing (MLST)

165 scheme alleles using mlst (https://github.com/tseemann/mlst), an in-house tool that uses a

166 BLAST algorithm [22] to search against the entire reference database of MLST profiles

167 (downloaded from https://pubmlst.org). In addition, acquired antimicrobial resistance genes,

168 including carbapenemases (blaKPC), were detected using another custom BLAST tool,

169 ABRicate (https://github.com/tseemann/abricate), to search against the ResFinder v2.1

170 database [21]

171

172 **Reconstructing 16S rRNA gene squences**

173 Near-complete 16S rRNA gene sequences were reconstructed from Patient B unassembled

174 short read metagenomic data (post removal of host-derived gDNA) using the Expectation

175 Maximization Iterative Reconstruction of Genes from the Environment (EMIRGE) program [23].

176 The following parameters were incorporated: the SILVA Small Subunit database was

177 employed as a training reference set, length of reads of 151, insert size of 683, standard

178 deviation of 68, and a phred score of 33 were selected to compute over 80 iterations.

179 Reconstructed 16S rRNA genes were queried against the Ribosomal Database Project using

180 BLAST. A k-mer based approach, using Kraken [24], classified unassembled read data to support

181 EMIRGE results.

182

183 **Reference genome assembly to validate metagenomic bins**

184 Reference genomes were assembled using Canu v1.5[25], trimmed

185 (https://github.com/tseemann/berokka) and circularized. Illumina short read data from the

186 same gDNA sample were used to correct and polish the draft PacBio genomes using Pilon

187 v1.22 [26] and Snippy v3.2 (https://github.com/tseemann/snippy). Further assembly of

188 unmapped short read data (i.e., Illumina reads that did not match chromosomal or larger

189 plasmid DNA from PacBio-derived data) using SPAdes v3.10.1 [27] ) was used to detect the

190 presence smaller plasmids potentially missed through DNA size selection. Prokka v1.11 [28] was

191 used to predict CDS regions and annotate the assembled genomes. Further characterisation

192 of reference genomes including multi-locus sequence typing and antimicrobial resistance

193 gene detection was performed in silico using the in-house developed tools, mlst and ABRicate

194 as described above.

195 **Transmission cluster inference**

196 To determine the most likely transmission cluster source for Patient A, previously sequenced

197 PacBio reference genomes from three local transmission clusters were assembled using the

198 methods described above, and used to build a custom Kraken database. The local

199 transmission clusters were defined through phylogenetic analysis of a maximum likelihood

200 tree described in Kwong *et al., (in prep)*. Whole-community metagenomic sequencing reads

201 were analysed in Kraken v0.10.5-beta [24] using the custom database to identify the most

202 closely related reference genome.

203

204

205 **Results**

206 *In silico* **typing and detection of AMR genes**

207 Analysis of the assembled metagenomic contigs from Patient A identified the presence of two

208 complete mlst profiles – ST258 *K. pneumoniae* and ST555 *E. faecium*. Table 1 highlights the

209 resistance AMR genes detected with 100% coverage. The gene encoding resistance to

210 carbapenem, *bla$_{KPC}$*, was detected at 100% coverage and nucleotide identity. The following

211 genes, with percentage coverage and nucleotide identity given in parentheses, were also

212 recovered from metagenomic data: *vanR-B* (100, 99.2), *vanS-B* (100, 99.6), *vanY-B* (100, 100),

213 *vanW-B* (100, 97.6), *vanH-B* (100, 99.4), *van-B* (100, 98.9), and *vanX-B* (100, 96.7), which

214 collectively is the vanB operon that encodes for vancomycin resistance in *Enterococcus*

215 *faecium*; while the vanB operon primarily encodes for antibiotic resistance in *E. faecium*,

216 Stinear et al., ([29]2001) have previously isolated vanB-positive anaerobic commensal bacteria

217 from human faeces. We hence reconstructed near-full length 16S rRNA genes to assign

218 taxonomy of operational taxonomic units in our metagenome.

219

220 **Taxonomy of metagenomic reads**

221 Near-full length 16S rRNA genes were reconstructed from the gut microbial community of

222 AUSMDU00008155 to detect the presence of *K. pneumoniae* and *E. faecium*. The Expectation

223 Maximization Iterative Reconstruction of Genes from the Environment program

224 reconstructed fifteen 16S rRNA genes, in which a *K. pneumoniae* 16S rRNA gene was

225 recovered with 100% nucleotide identity over 1161 bp (Table 2). Notably, given the presence

226 of the *vanB* operon, a near-complete 16S rRNA gene at 1344 bp was reconstructed and

227 classified as *E. faecium* at 100% nucleotide identity. Three uncultured organisms were

228 identified as belonging to *Bacteroitdetes*/*Bacteroides*, *Firmicutes*/*Clostridium* XIVa, and

229 *Proteobacteria*/*Sutterellaceae*, with 100% nucleotide identity, and over 1092 bp

230 reconstructed. An independent k-mer based approach (*i.e.,* Kraken) confirmed the presence

231 of *K. pneumoniae*, and *E. faecium* in AUSMDU00008155 metagenomic data.

232

233 **Reconstruction of isolate genomes from metagenomic reads**

234 Emergent Self Organizing Maps of the tetratnucleotide frequencies of AUSMDU00008155-

235 derived metagenomic contigs reconstructed discrete genome "bins" of isolates from the gut

236 microbial community (Figure 1A). Each point projected onto the ESOM represents DNA

9

237 fragments 2-5 kb in length, and colour coded with the following convention:

238 AUSMDU00008155 microbiome in red, AUSMDU00008118 in teal, AUSMDU00008119 in

239 navy, and an *E. faecium* AUS0085 strain in purple. A distinct *E. faecium* bin, and a largely

240 mixed bin consisting of closely related AUSMDU00008118 and AUSMDU00008119 derived

241 contigs were resolved (Figure 1B). Furthermore, a "satellite" cluster to the *K. pneumoniae* bin

242 consisted of only AUSMDU00008155 and AUSMDU00008119 contigs, and is typically

243 indicative of mobile genetic elements, such as, plasmids (Figure 1B, *circled*).

244

245 **Inference of transmission**

246 Using a custom Kraken database, we determined that Patient A's colonising ST258 *K.*

247 *pneumoniae* population were most closely related to the reference genome from

248 transmission cluster 2 (AUSMDU00008119), suggesting Patient A was most likely linked to this

249 transmission network (Figure 2). Of the local reference genomes (Kwong *et al., in prep*),

250 AUSMDU00008119 had 136 reads assigned, compared to 11 and 9 reads for the other cluster

251 references. This was corroborated by phylogenetic analysis of the multiple individual colony

252 sequences derived from Patient A's sample (Kwong *et al., in prep*). These data demonstrate

253 the potential of clinical metagenomics to guide source tracking and infection control efforts.

254

255

**Discussion**

Previous clinical metagenomic analyses have focused on understanding gut microbiota community diversity 16S rRNA gene level, or incorporating *in vitro* molecular diagnostics as the primary analysis to identify pathogenic isolates [12,30-32]. In contrast, in this proof-of-principle clinical study, we employed tetranucleotide frequency profiling as the main strategy to reconstruct genomes of community members directly from a faecal specimen, and determine the presence of KPC with strain-level resolution, as well as identifying previously unrecognized colonization with *vanB* VRE. Tetranucleotide profiles (i.e., frequencies of the 256 combinations of G,A,T, and C, in each contig) are a fundamental characteristic of DNA. Contigs with similar tetranucleotide profiles are derived from the same isolate, and therefore, projection of the frequency profiles of metagenomic contigs onto ESOMs can reconstruct isolate genomes independently of reference-based alignment approaches, such as BLAST and/or BWA MEM; this highlights the potential to allow strain level molecular characterization even when a reference isolate is not available *a priori* (c.f., [33-34]. The fact that this was achievable using a faeces sample further speaks to the validity of this approach, given the microbial complexity of this sample type.

Metagenomics permitted the comprehensive sampling of genomic content from an adult gut microbial community associated with KPC infection. We evaluated the AMR profile, and detected $bla_{KPC}$ at 100% coverage with zero gaps, and 100% nucleotide identity (Table 1) which encodes for carbapenem resistance in *Klebsiella pneumoniae*. Furthermore, the resolution of our analyses detected the presence of the genes encoding for an entire *vanB* operon, which confers vancomycin resistance in *E. faecium* isolates, at 100% coverage with zero gaps, and greater than 96% nucleotide identity (Table 1). Reconstruction of near full length 16S rRNA genes from metagenomic short read data recovered an 1161 bp 16S rRNA gene classified as *K. pneumoniae,* and 1344 bp 16S rRNA gene belonging to an *E. faecium* isolate, with 100% nucleotide identity and coverage (Table 2). Our metagenomic analyses uncovered VRE colonization in AUSMDU00008155, and initiated culture analysis of the faecal sample for the presence of VRE using selective media, and a retrospective report notifying the appropriate health care institution of a potential VRE carrier; a reporting that would otherwise have been missed by routine diagnostic analyses. Asymptomatic VRE carriage in AUSMDU00008155 may be facilitated by the co-occurrence of *C. bolteae* in the microbiome,

11

288    which is described to confer protection against VRE by Carballero *et al.,* ([30]2017); this

289    underscores the important roles microbial community members play in regulating infection

290    and immunity, and suggests whole community metagenomics could become a core

291    component of clinical microbiology.

292

293    Tetranucleotide frequency-ESOM analysis recovered metagenomic bins representing single

294    isolates from the AUSMDU00008155 gut microbiome (Figure 1A). Encouragingly, we found

295    that the number of reconstructed genomes (i.e., "bins") correlated with the number of 16S

296    rRNA genes recovered by our EMIRGE analysis. As a way of further validation, reference *E.*

297    *faecium* (Aus0085), and AUSMDU00008118 and AUSMDU00008119 genomes were included

298    to guide analysis of ESOMs. Figure 1B illustrates a distinct vanB-carrying *E. faecium* bin, and a

299    large *K. pneumoniae* bin with an associated "satellite" cluster. The main *K. pneumoniae* bin

300    consisted of both AUSMDU00008118 and AUSMDU00008119 derived contigs, while the

301    satellite cluster represents unique genomic content (most probably mobile genetic elements;

302    c.f., [16]) from AUSMDU00008119, and is therefore indicative of strain-level discrimination.

303    Although we have only assessed one faecal specimen in this study, the richness of microbial

304    characterization obtained from this sample using an alignment-free approach shows the

305    potential for applications in clinical and diagnostic microbiology. This potential is particularly

306    evident for clonal pathogens such as those examined here, where MLST results were unable

307    to distinguish between AUSMDU00008118, AUSMDU00008119, and metagenomic reads (i.e.,

308    all isolates were ST258).  For example, we accurately assigned Patient A's colonising

309    metagenomic *K. pneumoniae* isolate to a specific hospital infection cluster (Kwong *et al., in*

310    *prep*). Our ability to rapidly assign patient metagenomic isolates to pre-existing transmission

311    clusters identified during outbreak situations will better inform prevention control measures

312    to limit the spread of extensively drug-resistant pathogens across our hospital network.

313

314    In summary, the current study presents a clinical, and primarily culture-independent

315    investigation framework for the genomic profiling of patients colonised with multidrug-

316    resistant pathogens. Analysis of the whole community metagenome sampled from an adult

317    faecal sample revealed the presence of AMR genes conferring resistance to carbapenem and

318    vancomycin, and identification of pathogenic isolates. We have shown that metagenomic

319    binning, using tetranucleotide frequency profiles, can obtain strain-level resolution. A key

320    implication of incorporating metagenomics into routine clinical microbiology includes higher

321    resolution in AMR and pathogen detection, especially the detection of asymptomatic carriage

322    of antibiotic resistant microbes.

323

324    **Acknowledgements**

332

333

**References**

1. Williamson, D. A., Howden, B. P. & Stinear, T. P. *Mycobacterium chimaera* spread from heating and cooling units in heart surgery. *N Engl J Med* **376,** 600–602 (2017).

2. Zozaya-Vald s, E. *et al.* Target-specific assay for rapid and quantitative detection of *Mycobacterium chimaera* DNA. *J. Clin. Microbiol.* **55,** 1847–1856 (2017).

3. Phillips, A. *et al.* Whole genome sequencing of *Salmonella* Typhimurium illuminates distinct outbreaks caused by an endemic multi-locus variable number tandem repeat analysis type in Australia, 2014. *BMC Microbiology* 1–9 (2016). doi:10.1186/s12866-016-0831-3

4. Bakker, den, H. C. *et al.* Rapid Whole-Genome Sequencing for Surveillance of *Salmonella enterica* Serovar Enteritidis. *Emerg. Infect. Dis.* **20,** 1306–1314 (2014).

5. Staley, J. T. Measurement of *in situ* activities of nonphotosynthetic microorganisms in aquatic and terrestrial habitats. *Annual Review of Microbiology* **39,** 321–346 (1985).

6. Dreier, J. *et al.* Culture-negative infectious endocarditis caused by *Bartonella* spp.: 2 case reports and a review of the literature. *Diagnostic Microbiology and Infectious Disease* **61,** 476–483 (2008).

7. Richardson, D. C. *et al. Tropheryma whippelii* as a cause of afebrile culture-negative endocarditis: the evolving spectrum of Whipple's disease. *Journal of Infection* **47,** 170–173 (2003).

8. Nakamura, S. *et al.* Metagenomic Diagnosis of Bacterial Infections. *Emerg. Infect. Dis.* **14,** 1784–1786 (2008).

9. Pallen, M. J., Loman, N. J. & Penn, C. W. High-throughput sequencing and clinical microbiology: progress, opportunities and challenges. *Current Opinion in Microbiology* **13,** 625–631 (2010).

10. Miller, R. R., Montoya, V., Gardy, J. L., Patrick, D. M. & Tang, P. Metagenomics for pathogen detection in public health. *Genome Medicine* **5,** 1–114 (2013).

11. Mulcahy-O Grady, H. & Workentine, M. L. The Challenge and Potential of Metagenomics in the Clinic. *Front. Immunol.* **7,** 260 (2016).

12. Loman, N. J. *et al.* A Culture-Independent Sequence-Based Metagenomics Approach to the Investigation of an Outbreak of Shiga-Toxigenic *Escherichia coli* O104:H4. *JAMA* **309,** 1–9 (2013).

13. Lim, Y. W. *et al.* Clinical Insights from Metagenomic Analysis of Sputum Samples from Patients with Cystic Fibrosis. *J. Clin. Microbiol.* **52,** 425–437 (2014).

14. Freedman, A. J. E., Tan, B. & Thompson, J. R. Microbial potential for carbon and nutrient cycling in a geogenic supercritical carbon dioxide reservoir. *Environmental Microbiology* **31,** 533 (2017).

15. Wrighton, K. C. *et al.* Fermentation, hydrogen, and sulfur metabolism in multiple uncultivated bacterial phyla. *Science* **337,** 1661–1665 (2012).

16. Dick, G. J. *et al.* Community-wide analysis of microbial genome sequence signatures. *Genome Biology* **10,** R85 (2009).

17. Hug, L. A. *et al.* A new view of the tree of life. *Nature Microbiology* **1,** 1–6 (2016).

18. Carter, G. P. *et al.* Emergence of endemic MLST non-typeable vancomycin-resistant *Enterococcus faecium*. *Journal of Antimicrobial Chemotherapy* **71,**

381           3367–3371 (2016).

382    19.    Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for
383           Illumina sequence data. *Bioinformatics* **30,** 2114–2120 (2014).

384    20.    Peng, Y., Leung, C. M. H., Yui, S. M. & Chin, Y. L. F. IDBA-UD: A de Novo
385           Assembler for Single-Cell and Metagenomic Sequencing Data with Highly
386           Uneven Depth. *Bioinformatics* **11,** 1420–1428 (2012).

387    21.    Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes.
388           *Journal of Antimicrobial Chemotherapy* **67,** 2640–2644 (2012).

389    22.    Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. Basic Local
390           Alignment Search Tool. *Journal of Molecular biology* **215,** 403–410 (1990).

391    23.    Miller, C. S., Baker, B. J., Thomas, B. C., Singer, S. W. & Banfield, J. F. EMIRGE:
392           reconstruction of full-length ribosomalgenes from microbial community short
393           readsequencing data. *Genome Biology* **12,** 1–14 (2011).

394    24.    Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence
395           classification using exact alignments. *Genome Biology* **15,** 1–12 (2014).

396    25.    Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-
397           mer weighting and repeat separation. *Genome Research* **27,** 722–736 (2017).

398    26.    Walker, B. J. *et al.* Pilon: An Integrated Tool for Comprehensive Microbial
399           Variant Detection and Genome Assembly Improvement. *PLoS ONE* **9,** e112963
400           (2014).

401    27.    Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its
402           Applications to Single-Cell Sequencing. *Journal of Computational Biology* **19,**
403           455–477 (2012).

404    28.    Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30,**
405           2068–2069 (2014).

406    29.    Stinear, T. P., Olden, D. C., JOhnson, P. D. R., Davies, J. K. & Grayson, M. L.
407           Enterococcal *vanB* resistance locus in anaerobic bacteria in human faeces. *The*
408           *lancet* **357,** 855–856 (2001).

409    30.    Caballero, S. *et al.* Cooperating Commensals Restore Colonization Resistance to
410           Vancomycin-Resistant *Enterococcus faecium*. *Cell Host and Microbe* **21,** 592–
411           602.e4 (2017).

412    31.    Zhou, Y. *et al.* Metagenomic Approach for Identification of the Pathogens
413           Associated with Diarrhea in Stool Specimens. *J. Clin. Microbiol.* **54,** 368–375
414           (2016).

415    32.    Yu, H.-J. *et al.* International Journal of Infectious Diseases. *International Journal*
416           *of Infectious Diseases* **53,** 30–33 (2016).

417    33.    Sharon, I. *et al.* Time series community genomics analysis reveals rapid shifts in
418           bacterial species, strains, and phage during infant gut colonization. *Genome*
419           *Research* **23,** 111–120 (2013).

420    34.    Morowitz, M. J. *et al.* Strain-resolved community genomic analysis of gut
421           microbial colonization in a premature infant. *Proceedings of the National*
422           *Academy of Sciences* **108,** 1128–1133 (2011).

423

424

425

426

427 **Figures and Tables**

428 **Table 1: Detection of antimicrobial resistance genes in AUSMDU00008155 metagenome**

429 **using ABRicate**

| Contig | Start | End | Gene | Gaps | % Coverage | % Identity | Accession |
|---|---|---|---|---|---|---|---|
| FM_4827# | 2150 | 3076 | *sul1_2* | 0 | 100 | 100 | CP002151 |
| FM_1035 | 10808 | 11689 | *blaAUSMDU0 0008119_1* | 0 | 100 | 100 | AY034847 |
| FM_1074 | 7275 | 8135 | *blaSHV-12_1* | 0 | 100 | 100 | AF462395 |
| FM_10044 | 186 | 1001 | *aph(3')-Ia_1* | 0 | 100 | 100 | V00359 |
| FM_452 | 22060 | 22866 | *VanY-B_1* | 0 | 100 | 100 | AF192329 |
| FM_784 | 15763 | 16260 | *dfrG_1* | 0 | 100 | 100 | AB205645 |
| FM_313 | 6894 | 10046 | *oqxB_1* | 0 | 100 | 100 | EU370913 |
| FM_313 | 10070 | 11245 | *oqxA_1* | 0 | 100 | 100 | EU370913 |
| FM_2352 | 2933 | 4853 | *tet(W)_4* | 0 | 100 | 99.9 | AJ427422 |
| FM_10870 | 34 | 873 | *blaOXA-9_2* | 0 | 100 | 99.9 | JF703130 |
| FM_9051 | 161 | 1381 | *tet(40)_1* | 0 | 100 | 99.8 | FJ158002 |
| FM_7453 | 1068 | 1565 | *dfrA12_1* | 0 | 100 | 99.8 | AB571791 |
| FM_12072 | 99 | 965 | *aadE_1* | 0 | 100 | 99.8 | KF864551 |
| FM_452 | 20546 | 21889 | *VanS-B_1* | 0 | 100 | 99.6 | AF192329 |
| FM_452 | 23708 | 24679 | *VanH-B_1* | 0 | 100 | 99.4 | AF192329 |
| FM_452 | 19884 | 20546 | *VanR-B_1* | 0 | 100 | 99.2 | AF192329 |
| FM_452 | 24672 | 25700 | *VanA-B_1* | 0 | 100 | 98.9 | AF192329 |
| FM_137 | 6369 | 7847 | *msr(C)_1* | 0 | 100 | 98.9 | AY004350 |
| FM_255 | 25338 | 25757 | *fosA_3* | 0 | 100 | 98.6 | NZ_ACWO01000079 |
| FM_452 | 22884 | 23711 | *VanW-B_1* | 0 | 100 | 97.6 | AF192329 |
| FM_452 | 25706 | 26314 | *VanX-B_1* | 0 | 100 | 96.7 | AF192329 |

430 #FM_xxxxx = AUSMDU00008155 and the associated metagenomic contig number

431

432

433

434

435

436

437

16

438 **Table 2: Reconstructed 16S rRNA genes from AUSMDU00008155 metagenome using**

439 **EMIRGE**

| Taxon | Subject length (bp) | Query length (bp) | Nucleotide ID | % Identity |
|---|---|---|---|---|
| *Akkermansia Muciniphila* | 1434 | 1500 | 1431 | 99.8 |
| *Bacteroides cellulosilyticus* | 1193 | 1202 | 1183 | 99.2 |
| *Bacteroides uniformis* | 1145 | 1145 | 1143 | 99.8 |
| *Bacteroides uniformis* | 680 | 680 | 676 | 99.4 |
| *Clostridium bolteae* | 1309 | 1309 | 1305 | 99.7 |
| *Clostridium glycyrrhizinilyticum* | 1268 | 1266 | 1231 | 97.1 |
| *Dialister invisus* | 1290 | 1302 | 1287 | 99.8 |
| *Enterococcus faecium* | 1344 | 1344 | 1344 | 100 |
| *Eubacterium dolichum* | 1484 | 1516 | 1381 | 93.1 |
| *Klebsiella pneumoniae* | 1161 | 1161 | 1161 | 100 |
| *Lactobacillus pentosus* | 1468 | 1468 | 1467 | 99.9 |
| *Parabacteroides merdae* | 1376 | 1376 | 1376 | 100 |
| Uncultured bacterium[a] | 1092 | 1092 | 1092 | 100 |
| Uncultured organism[b] | 1121 | 1121 | 1121 | 100 |
| Uncultured organism[c] | 1191 | 1191 | 1191 | 100 |

440 [a]*Bacteroidetes/Bacteroides*
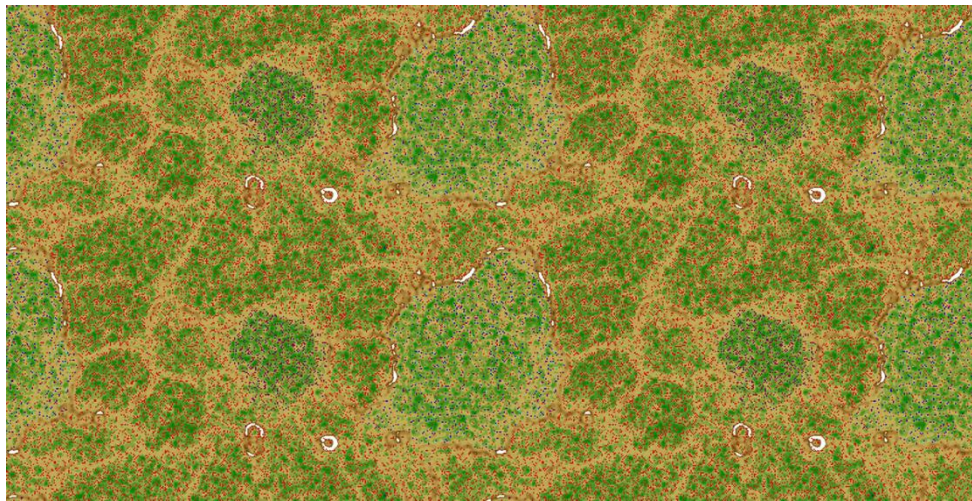
441 [b]*Firmicutes/Clostridium* XIVa
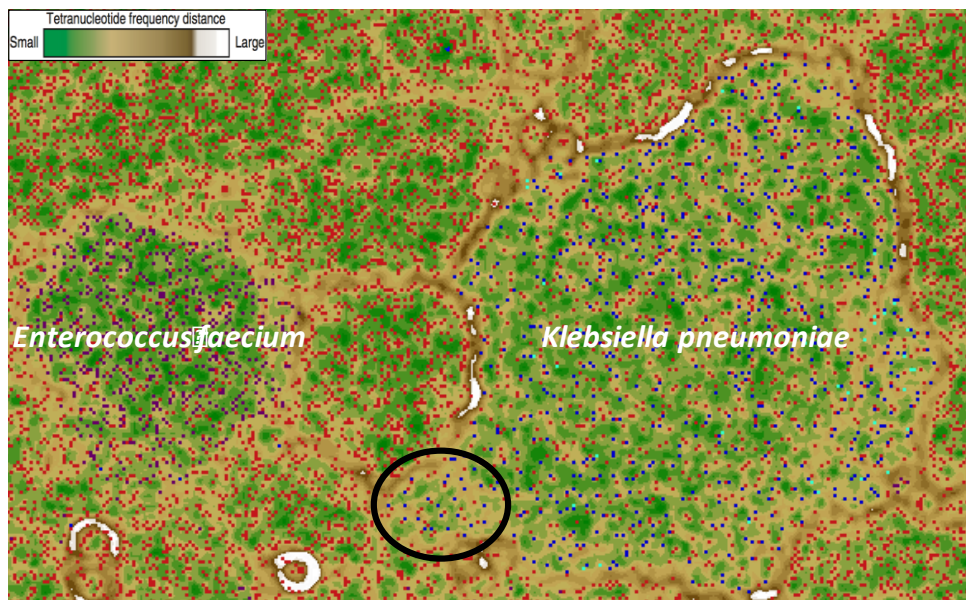
442 [c]*Proteobacteria/Sutterellaceae*

443

444 **Figures:**

445 **Figure. 1**

446 **A**



455 **B**



**Figure 1. (A) Emergent Self Organizing Map of the tetranucleotide frequencies of AUSMDU00008155 contigs in binned genomes representing (B) *Klebsiella pneumoniae* and *Enterococcus faecium*.** Tetranucleotide frequency profiles of DNA fragments between 2 kb to 5 kb in length are projected onto the ESOM. A green background indicates small tetranucleotide frequency distances, white background represents large tetranucleotide frequency distances, while contigs derived from AUSMDU00008155 microbiome is highlighted in red, AUSMDU00008118 in teal, AUSMDU00008119 in navy, and a *vanB*-carrying *E. faecium* in purple.
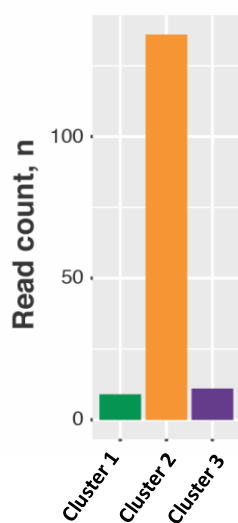
18

474 **Fig. 2**

475
476
477
478
479
480
481
482
483
484
485



486 **Figure 2: Metagenomic attribution of transmission clusters.** The bar graph shows the
487 number of metagenomic reads assigned to each of the reference genomes, representing
488 different transmission clusters described in Kwong *et al.,* (*in prep*).