1 **Generating closed bacterial genomes from long-read nanopore sequencing of**

2 **microbiomes**

3 Eli L. Moss[1], Ami S. Bhatt[1,**]

4 [2] Department of Medicine (Hematology, Blood and Marrow Transplantation) and Department of

5 Genetics, Stanford University, Stanford, California, USA.

6 [**] To whom correspondence should be addressed: asbhatt@stanford.edu

7 **Abstract**

8        We present the first method for efficient recovery of complete, closed genomes directly

9    from microbiomes using nanopore long-read sequencing and assembly. We apply our approach

10    to three healthy human gut communities and compare results to short read and read cloud

11    approaches. We obtain nine finished genomes including the first reported closed genome of

12    *Prevotella copri*, an organism with highly repetitive genome structure prevalent in non-western

13    human gut microbiomes.

14 **Main Text**

15      *De novo* reconstruction of complete microbial genomes from metagenomes has been a

16 longstanding goal of microbiome research. Although current reference-based methods are able

17 to detect known organisms and genes in metagenomes, only *de novo* approaches are able to

18 characterize novel genome sequences, or accurately place mobile or transferred elements in

19 new genomic contexts. The tremendous diversity and plasticity of bacterial genomes, as well as

20 the difficulty of bacterial isolation and culture, demand effective culture-free methods for

21 producing genomes directly from metagenomes.

22      Current metagenomic sequencing and assembly methods do not typically yield finished

23 bacterial genomes, although previous efforts have achieved single closed genomes in simple

24 communities[1], or multiple genomes with skilled manual assembly and scaffolding[2].

25 Consequently, genome drafts are formed by grouping (i.e. binning) similar contigs within

26 fragmentary assemblies[3,4]. This is an imperfect process, often compromising the purity or the

27 completeness of the genome reconstruction. As assembly contiguity increases, the sensitivity

28 and specificity of genome binning are improved as fewer, larger contigs need to be grouped to

29 form each genome. Indeed, at the point when genomes are assembled in single contigs, binning

30 becomes unnecessary. Nanopore long read assembly has yielded complete genomes in

31 cultured bacterial isolates[5–8], suggesting potential for effective assembly in more complex

32 microbial communities. However, the performance of nanopore and other long read approaches

33 in metagenomic sequencing and assembly has been limited by the lack of effective and efficient

34 methods to maximize molecular weight, mass yield and purity of DNA extracted from these

35 samples.

36      We present a workflow consisting of stool DNA extraction, nanopore sequencing,

37 assembly and post-processing steps capable of producing multiple complete, circular bacterial

38 genomes directly from metagenomes. Our extraction approach produces microgram quantities

39 of pure, high molecular weight (HMW) DNA suitable for long read sequencing from as little as

40  300 milligrams (mg) of stool. Our computational workflow, consisting of assembly and post-

41  processing, does not involve manual intervention in assembly, scaffolding, bacterial isolation, or

42  existing reference coverage of the target metagenome. Thus, this workflow is the first to provide

43  a rapid, simple, cost-effective, automated approach to close high numbers of bacterial genomes

44  directly from metagenomic samples.

45  Short read and read cloud data and assemblies for samples P1 and P2-A were used

46  without modification as previously described[9]. The standard approach used to extract DNA for

47  these libraries produced fragmented DNA which incurred a severe loss with size selection,

48  necessitating approximately 300 mg input stool to assure the 1 nanogram (ng) final HMW DNA

49  mass required for read cloud library preparation. Current long read library prep protocols require

50  1000 ng of HMW input DNA, well beyond the practical capability of existing stool DNA extraction

51  techniques. In order to maximize the throughput and read length of nanopore sequencing, a

52  new approach yielding DNA in dramatically higher quantity and molecular weight was needed.

53  We developed a method for HMW extraction capable of yielding 1000-fold more DNA

54  over 5 kb than a conventional bead-beating approach (Supplementary Figure 1, see Methods).

55  We applied this method to two samples (P1 and P2-A) as well as a third sample (P2-B),

56  collected 15 months later from the second individual. HMW DNA extraction yielded at least 1 µg

57  HMW DNA per 300 mg input stool mass for all samples (Supplementary Table 1). Nanodrop

58  measurement produced $A_{260/280}$ ratios over 1.86 and $A_{260/230}$ ratios over 2.23 for all samples,

59  indicating absence of contaminants such as proteins, solvents and salts.

60  We obtained a total of 12.7 giga-base pairs (Gbp), 6.1 Gbp, and 7.6 Gbp of long read

61  data for samples P1, P2-A, and P2-B, respectively (Supplementary Figure 2, Supplementary

62  Table 2) with N50 values of 4.7 kbp, 3 kbp and 3 kbp. The taxonomic composition of reads

63  obtained through our approach was compared to that obtained by standard mechanical lysis

64  and short read sequencing methods (see Methods). Although precise rank order relative

65  abundances varied, we noted higher Shannon diversity from the present approach (P2: 2.0 vs.

66    1.14; P1: 2.0 vs. 1.8). We also detected all genera represented by more than 200 short reads

67    from the traditional short read sequencing in the long read data.

68         Our assembly and post-processing workflow yielded whole-assembly N50 values of 453

69    kbp, 571 kbp and 564 kbp for the three samples P1, P2-A and P2-B. In comparison, the short

70    read approach did not exceed assembly N50 of 34 kbp across samples P1 and P2-A, in spite of

71    3- to 6-fold more read data (37-38 Gbp). Our approach also surpassed the read cloud N50

72    values of 116 kbp and 12 kbp. However, read cloud and short read assemblies were between

73    1.5- and 2.1-fold larger than corresponding long read assemblies, likely due to the much greater

74    volume of raw data available from these datasets (Supplementary Table 3).

75         Contigs from each approach were binned to form draft genomes, which were evaluated

76    and assigned 'High Quality', 'Complete' and 'Incomplete' labels as described[9]. Briefly, drafts at

77    least 90% complete and with at most 5% contamination are termed 'Complete', and drafts also

78    containing at least one each of the 5S, 16S and 23S rRNA loci, as well as at least 18 tRNA loci,

79    are labeled 'High Quality'. All others are 'Incomplete'. While read cloud and short read methods

80    produced more complete bins and a comparable number of high quality bins compared to the

81    long read approach, the long read approach produced bins with much higher contiguity (Figure

82    1). The present approach yielded nine high quality genomes with N50 over 2 Mbp, whereas the

83    read cloud approach yielded only one. Short read bins never exceeded 550 kbp. Finally, the

84    present approach yields a comparable quantity of high quality genomes at far higher contiguity

85    with lower capital equipment requirement, sequencing cost and turnaround time (Supplementary

86    Table 4).

87         Nanopore long read assembly yielded nine complete, circularizeable bacterial genomes

88    across the three sequenced samples, and a maximum of four from a single sample (P1),

89    compared to zero from the short read, read cloud, and synthetic long read approaches

90    previously applied to these samples[9]. Assembled genomes are up to 5Mbp in length, and in

91    several cases (*Prevotella copri, Subdoligranulum variabile, Phascolarctobacterium faecium,* and

92    *Bacteroides uniformis*) represent the first closed genomes for their species. Closed genomes

93    ranged in coverage depth between 75 (*Oscillibacter sp.)* and 785 (*P. copri)*. Closed genomes

94    were largely structurally concordant and similar in sequence to existing published genome

95    sequences (Supplementary Figure 3, Supplementary Table 5), although in some cases we do

96    note extensive strain divergence; for example, our closed *Dialister invisus* genome exhibits

97    multiple large-scale inversions relative to the available reference (see below).

98         Completed bacterial genomes included two for *Prevotella copri* in samples P2-A and P2-

99    B. This organism lacks a closed reference, in spite of extensive efforts to assemble *P. copri* and

100   other members of the genus *Prevotella*[10]. Our previous efforts using read clouds, short reads

101   and synthetic long reads to assemble these communities also had limited success with this

102   organism, never exceeding a genome N50 of 130 kbp, in spite of coverage depth in excess of

103   4,800x[9]. The two *P. copri* genomes obtained from samples separated by 15 months display high

104   concordance, with 99.94% of bases aligned and 99.89% nucleotide identity, suggesting nearly

105   identical strain composition in the two time points.

106         The difficulty of assembling the *P. copri* genome stems from its high degree of sequence

107   repetition. A direct assembly of highly abundant *k*-mers (*k*=101, occurring more than 5 times)

108   found in our complete genome assembly yielded two insertion sequences (ISs) (see Methods),

109   one 1.1 kbp IS66 family sequence and one 1.6 kbp IS1380 family sequence. These were found

110   to be assembled in a total of 29 genomic loci between the two timepoints, but IS instances

111   absent from the consensus assembly were detected directly in long reads at an additional 45

112   loci (see Methods). These insertion sites, whether fixed in the strain population or varying

113   between strains, co-locate with breaks in short read and read cloud assemblies, illustrating their

114   impact on these types of assembly (Figure 2).

115         Other complete genomes include *Phascolarctobacterium faecium* assembled in samples

116   P2-A and P2-B at relative abundances of 4% and 1.41%, respectively. These assemblies are

117   the first complete genomes for this species, and display high structural and nucleotide

118    concordance with the closest available reference (Supplementary Figure 3; 98.9% identity,

119    88.5% sequence alignment) and with each other (99.81% identity, 99.97% sequence

120    alignment). Sample P2-A also yielded the first circular genome for *Dialister invisus*, present in

121    that sample at 1.03% relative abundance. We find similar structural divergence compared to the

122    available reference (Supplementary Figure 3), and concordance with the read cloud draft, which

123    contained identical large-scale structural inversions (99.90% identity, 99.97% sequence

124    alignment) (Figure 2, Supplementary Table 5). Although the *Dialister invisus* assembled in

125    sample P2-B was assessed complete, it was not found to be circularizable.

126          In sample P1, we obtained circular genomes for *Bacteroides uniformis* (6% abundance

127    in long reads)*, Alistipes finegoldii* (2% abundance)*, Oscillibacter sp.* (0.14% abundance)*, and

128    *Subdoligranulum variabile* (0.37% abundance). Of these, we were able to obtain structurally

129    concordant reference sequences for all but *Subdoligranulum*, for which we could not locate a

130    reference with more than 19% aligned bases, suggesting the possibility of a novel strain. Seven

131    16S rRNA loci were assembled in this genome, all bearing 98% sequence identity to the closest

132    match from *Subdoligranulum variabile* strain BI114, for which no genome reference is available

133    for comparison. Identity with read cloud assemblies in all cases was over 99.7%, with over

134    98.3% of bases aligning to the assembled draft (Supplementary Table 5). For all closed

135    genomes, read cloud and short read assembly yielded more fragmentary assemblies, which

136    were only partially recovered by binning (Figure 2).

137          Our approach relies on consensus refinement based on short read data to correct

138    homopolymer errors intrinsic to the current nanopore sequencing technology. Although long

139    read-based consensus refinement is possible and partially effective, we find that it cannot fully

140    replace short read correction (Supplementary Figure 4). We found that uncorrected long read

141    assembly demonstrated a 3% error rate with 3-mer homopolymers, assembled too short by an

142    average of 0.5 nucleotides. This worsens to a 65% error rate on 6-mer homopolymers, which

143    were assembled too short by an average of 1.3 nucleotides. On average, 63 homopolymers of

144    length 3 or greater were found per kilobase of assembled sequence, of which 4.5 (7.1%) were

145    found to require correction with short reads. CheckM, a tool which annotates genome

146    completeness based on single copy core gene detection, demonstrates a low detection rate on

147    uncorrected assemblies, consequently under-reporting genome completeness even on

148    circularizeable whole-genome contigs. For instance, the genome annotated *Oscillibacter sp.* in

149    sample P1 is annotated 38% complete in the uncorrected assembly. This rises to 68% complete

150    after correction with long reads. With short read correction, the genome receives a 96%

151    completeness annotation, compared to 98% for the sole available closed genome reference

152    sequence (strain PEA192). The present workflow for sequencing and assembly can operate

153    solely with long reads and will yield structurally correct and complete genomes, although with

154    reduced nucleotide accuracy. Future advances in nanopore sequencing technology that

155    decrease the homopolymer-repeat related errors will likely lessen or remove the requirement for

156    supplemental short read sequencing to achieve genomes with high nucleotide fidelity.

157         Although the present approach has achieved effective assembly of bacterial genomes

158    from metagenomes, we anticipate that future advances in metagenomic DNA extraction

159    methods and nanopore long read assembly will improve read length and reduce the read

160    coverage required to close genomes. In addition, epigenetic modification detection will add to

161    future metagenomic studies by revealing phage and bacterial sequence methylation patterns,

162    methylation-based contig binning approaches, and epigenetic regulation of bacterial DNA-

163    protein interactions.

164         In conclusion, our approach assembles the first complete genome of *Prevotella copri*, an

165    organism with high prevalence in non-western guts and with emerging, potentially strain-specific

166    links to human health and disease[11,12]. The high copy number of IS66 and IS1380 family

167    insertion sequences in this genome limit the effectiveness of short read approaches, despite

168    receiving over 4,800x coverage in an earlier metagenomic sequencing study[9] and extensive

169    isolate sequencing in a separate effort[10]. IS1380 has been previously reported to carry an

170     outward-facing promoter capable of upregulating adjacent gene sequences, and has been

171     found to impact antibiotic resistance gene regulation and resistance phenotype[13]. We anticipate

172     that this approach will help illuminate the role of repetitive classes of genomic elements with

173     important effects on cellular and clinical phenotypes, and facilitate efforts to broaden human

174     microbiome research to global populations where *Prevotella* are highly prevalent[14]. Closing

175     these and other genomes will allow investigation into the complete functional repertoire and

176     potential phenotypes of individual microbes, even when these organisms are difficult to culture

177     or are found in mixed communities, facilitating future research in important human microbiomes

178     and poorly characterized microbial communities such as soil and marine sediment.

179  **Bibliography**

180  1.  Leonard, M.T. et al. *Front. Microbiol.* **5**, 361 (2014).

181  2.  Anantharaman, K. et al. *Nat. Commun.* **7**, 13219 (2016).

182  3.  Bowers, R.M. et al. *Nat. Biotechnol.* **35**, 725–731 (2017).

183  4.  Kang, D.D., Froula, J., Egan, R. & Wang, Z. *PeerJ* **3**, e1165 (2015).

184  5.  Koren, S. & Phillippy, A.M. *Curr. Opin. Microbiol.* **23**, 110–120 (2015).

185  6.  Risse, J. et al. *Gigascience* **4**, 60 (2015).

186  7.  Turner, D.J., Dai, X., Mayes, S. & Juul, S. *bioRxiv* 026930 (2015).

187  8.  Loman, N.J., Quick, J. & Simpson, J.T. *Nat. Methods* **12**, 733–735 (2015).

188  9.  Bishara, A. et al. *Nat. Biotechnol.* (2018).doi:10.1038/nbt.4266

189  10. Gupta, V.K., Chaudhari, N.M., Iskepalli, S. & Dutta, C. *BMC Genomics* **16**, 153 (2015).

190  11. Scher, J.U. et al. *eLife Sciences* **2**, e01202 (2013).

191  12. Pianta, A. et al. *Arthritis & Rheumatology* **69**, 964–975 (2017).

192  13. Kato, N., Yamazoe, K., Han, C.-G. & Ohtsubo, E. *Antimicrob. Agents Chemother.* **47**, 979–

193      985 (2003).

194  14. Schnorr, S.L. et al. *Nat. Commun.* **5**, 3654 (2014).

195  15. Mukhopadhyay, T. & Roth, J.A. *Nucleic Acids Res.* **21**, 781–782 (1993).

196  16. Not provided, R. & Schwessinger, B.doi:10.17504/protocols.io.n7hdhj6

197  17. Koren, S. et al. *Genome Res.* (2017).doi:10.1101/gr.215087.116

198  18. Chakraborty, M., Baldwin-Brown, J.G., Long, A.D. & Emerson, J.J. *Nucleic Acids Res.* **44**,

199      e147 (2016).

200  19. Hunt, M. et al. *Genome Biol.* **16**, 294 (2015).

201  20. at <https://nanoporetech.github.io/medaka/index.html>

202  21. Walker, B.J. et al. *PLoS One* **9**, e112963 (2014).

203  22. Danecek, P., McCarthy, S., Li, H. & Others (2015).

204    23. Delcher, A.L., Salzberg, S.L. & Phillippy, A.M. *Curr. Protoc. Bioinformatics* **00**, 10.3.1–

205         10.3.18 (2003).

206    24. Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P. & Tyson, G.W. *Genome Res.*

207         **25**, 1043–1055 (2015).

208    25. Wood, D. & Salzberg, S. *Genome Biol.* **15**, R46 (2014).

209    26. Dixon, P. *J. Veg. Sci.* **14**, 927–930 (2003).

210    27. Wickham, H. (Springer Science & Business Media: 2009).

211    28. Hahne, F. & Ivanek, R. *Methods Mol. Biol.* **1418**, 335–351 (2016).

212    29. Højsgaard, S., Højsgaard, M.S. & Hmisc, D. *The Newsletter of the R Project* **6**, 1 (2006).

213    30. Wickham, H. *R package version* **1**, (2012).

214    31. Köster, J. & Rahmann, S. *OASIcs-OpenAccess Series in Informatics* **26**, (2012).

215    32. Marçais, G. & Kingsford, C. *Bioinformatics* **27**, 764–770 (2011).

216    33. Bankevich, A. et al. *J. Comput. Biol.* **19**, 455–477 (2012-5).

217    34. Li, H. *Bioinformatics* **34**, 3094–3100 (2018).

218 **Methods**

219 DNA extraction

220       Short read and read cloud libraries were prepared as previously described[9]. Previously,

221 DNA was extracted from samples P1 and P2-A with a commercial extraction kit using bead-

222 beating lysis.

223       For high molecular weight (HMW) extraction, approximately 0.7g frozen stool was

224 aliquoted into 2mL eppendorf tubes (Eppendorf, Hamburg, Germany) with a 4mm biopsy punch

225 (Integra Miltex, Plainsboro, NJ) and suspended in 500μL PBS (Fisher Scientific, Waltham, MA)

226 with brief gentle vortexing. 5uL of lytic enzyme solution (Qiagen, Hilden, Germany) was added

227 and the samples were mixed by gentle inversion six times, then incubated for one hour at 37°C.

228 12μL 20% (w/v) SDS (Fisher Scientific, Waltham, MA) was added with approximately 100μL

229 vacuum grease (Dow-Corning, Midland, MI) functioning as phase lock gel[15]. 500μL phenol

230 chloroform isoamyl alcohol at pH 8 (Fisher Scientific, Waltham, MA) was added, and samples

231 were gently vortexed for five seconds, then centrifuged at 10,000g for five minutes with Legend

232 Micro 21 microcentrifuge (Fisher Scientific, Waltham, MA). The aqueous phase was then

233 decanted into a new 2mL tube.

234       Next, DNA was precipitated with 90μL 3M sodium acetate (Fisher Scientific) and 500uL

235 isopropanol (Fisher Scientific) for ten minutes at room temperature. After inverting three times

236 slowly, samples were incubated at room temperature for 10 minutes, then centrifuged 10

237 minutes at 10,000g. The supernatant was removed and the pellet was washed two times with

238 freshly prepared 80% (v/v) ethanol (Fisher Scientific). The pellet was then air dried with heating

239 for ten minutes at 37°C or until the pellet was matte in appearance, and then resuspended in

240 100μL nuclease-free water (Ambion, Thermo Fisher Scientific, Waltham, MA). 1mL Qiagen

241 buffer G2, 4μL Qiagen RNase A at 100mg/mL, and 25μL Qiagen Proteinase K were added, the

242 samples were then gently inverted three times, and then were incubated 90 minutes at 56°C.

243 After the first 30 minutes, pellets were dislodged by a single gentle inversion.

244     One Qiagen Genomic-tip 20/G column per sample was equilibrated with 1mL Qiagen

245     buffer QBT and allowed to empty by gravity flow.  Samples were gently inverted twice, applied

246     to columns and allowed to flow through.  Three stool extractions were combined per column.

247     Columns were then washed with 3mL Qiagen buffer QC, then DNA was eluted with 1mL Qiagen

248     buffer QF prewarmed to 56°C.  Eluted DNA was then precipitated by addition of 700µL

249     isopropanol followed by inversion and centrifugation for 15 minutes at 10,000g.  The

250     supernatant was carefully removed by pipette, and pellets were washed with 1mL 80% (v/v)

251     ethanol.  Residual ethanol was removed by air drying ten minutes at 37°C. This was followed by

252     resuspension of the pellet in 100µL water overnight at 4°C without agitation or any kind.

253     DNA was then size selected with a modified SPRI bead protocol as described [16], with

254     minor modifications: beads were added at 0.9x, and eluted DNA was resuspended in 50µL

255     water.  The concentration, purity and fragment size distribution of extracted DNA was then

256     quantified with the Qubit fluorometer (Thermo Fisher Scientific, Waltham, MA), Nanodrop

257     (Thermo Fisher Scientific), and Tapestation 2200 (Agilent Technologies, Santa Clara, CA),

258     respectively (Supplementary Table 1).

259

260     <u>Sequencing</u>

261     Extracted DNA samples were prepared for long read sequencing with the Oxford

262     Nanopore Technologies (ONT, Oxford, UK) Ligation library preparation kit according to the

263     manufacturer's standard protocol.  Libraries were sequenced with the ONT MinION sequencer

264     using rev C R9.4 flow cells, allocating one flowcell per sample. The sequencer was controlled

265     with the MinKNOW v2.2.12 software running on a MacBook Pro (model A1502, Apple,

266     Cupertino, CA), with data stored to a Vectotech 2Tb SSD hard drive.  Sequencing runs were

267     scheduled for 48 hours, and allowed to run until fewer than 10 pores remained functional. After

268     sequencing, data were uploaded to the Stanford Center for Genomics computational cluster for

269      analysis (see below). Short read libraries were prepared and sequenced as described

270      previously[9].

271

272      <u>Sequence assembly and analysis</u>

273      Raw data were basecalled with Albacore v2.3.1, and assembled in two separate runs

274      with Canu v1.7.1 with the -nanopore preset parameter [17]. The two runs differed by the estimated

275      genomeSize parameter, provided as either 50m or 100m. The two separate assemblies were

276      then merged with quickmerge v0.40[18], circularized with Circlator v1.5.5[19] and Encircle (present

277      study, see below), and then polished with either Medaka [20] or a parallelized version of Pilon

278      v1.22 [21] (present study) for long read or short read consensus refinement, respectively. In order

279      to parallelize Pilon, reference sequences were divided into 100kb segments, short reads

280      aligning to each segment were downsampled to at most 40x coverage depth, and Pilon was

281      used to detect errors within the reference and read subset. These errors were then aggregated

282      across all subset runs and used to generate a refined consensus with bcftools[22]. Errors found in

283      homopolymers were identified with an in-house script, homopolymer_error_analyzer.

284      Sequences are binned and annotated as previously described [9].

285      There is presently no straightforward, comprehensive method for determining circularity

286      in metagenome-assembled genomes. A minority of the circular genomes we obtained (*D.*

287      *invisus* in P2-A, *P. copri and Phascolarctobacterium* in P2-B) were circularized by an existing

288      genome circularization tool. In several cases, assembled genome contigs extended beyond the

289      wrap-around point of the circular chromosome, resulting in what we term over-circularization

290      (supp fig 5). Over-circularized contigs contain redundant sequences at their termini which

291      spuriously increase apparent contamination when assessed by CheckM. In order to trim over-

292      circularized contig ends in order to obtain a nonredundant, circular genome, we developed

293      Encircle, a utility which performs contig self-alignment with Mummer[23] and determines when

294      over-circularization has taken place, then outputs precise trim coordinates to circularize the

295    genome.  The genomes of *P. copri, Phascolarctobacterium sp.,* and *Dialister invisus* in sample

296    P2-A, as well as *Oscillibacter sp.* and *Subdoligranulum* (supp figure 5) in sample P1, were over-

297    circularized and required trimming.  In addition, the genomes of *Bacteroides uniformis* and

298    *Alistipes finegoldii* were determined to be circular by concatenating the first and last 20kbp of

299    the assembled genome, mapping long reads to the junction, and inspecting alignments for

300    reads spanning the gap; *B. uniformis* was found to be slightly overcircularized by 10kbp (below

301    the limit of detection of Encircle), and *A. finegoldii* was found to be perfectly circularizeable.

302        Binning was performed and evaluated as previously described[9]. Due to the complete

303    genomes present in our assemblies, binning became unnecessary for some organisms, and

304    instead led to several cases of genomic contamination as assessed with CheckM.  In cases

305    where >5% contamination occurred in a bin with one genome-scale contig and several much

306    smaller (<100kbp) sequences, the smaller sequences were removed and the largest sequence

307    was re-evaluated with CheckM[24], in two cases yielding complete and uncontaminated genomes.

308        Long and short reads were taxonomically classified with Kraken[25], and Shannon

309    diversity was calculated with vegan[26]. Figures were generated with ggplot2[27], gviz[28], doBy[29] and

310    reshape2[30].  All workflows were implemented with Snakemake[31].

311

312    Insertion sequence strain diversity

313        K-mers represented more than 6 times in the *Prevotella copri* assemblies were identified

314    with Jellyfish2[32].  These were assembled with SPAdes[33] two obtain two full-length insertion

315    sequences. These sequences were located in the genome assemblies by alignment with

316    minimap2[34].  In order to locate additional unassembled insertion sequences present in strains of

317    *P. copri*, reads containing insertion sequences were identified by alignment with minimap2, then

318    200 bases immediately upstream of the insertion sequence were taken from each read and

319    aligned to the genome assembly.

320    In order to quantify the relative abundance of *P. copri* strains carrying each IS instance,

321    long reads were first aligned to the assembled IS sequences. Long reads containing IS

322    sequences were isolated, and flanking sequences 200bp upstream of the IS were extracted and

323    realigned to the genome assembly. IS-flanking sequence depth was compared to local overall

324    coverage depth to obtain the relative abundance of strains carrying a given IS. Only 18 insertion

325    sites carried fixed ISs and a further 56 sites showed a mixture of strains with and without an IS

326    (Figure 2).

327

328    Data availability

329    All sequence data, whole metagenome assemblies and individual completed genomes

330    can be found at NCBI BioProject under accession PRJNA508395.

331

332    Code availability

333    All workflows and associated environments and tools can be found at

334    https://github.com/elimoss/metagenomics_workflows/.

342   **Competing financial interests**

343   The authors declare no competing financial interests.

344

345   **Figure Legends**

346   Figure 1

347   Taxonomic read composition and per-organism assembly contiguity for healthy gut assemblies,

348   overall genome draft counts in two healthy human gut microbiomes (samples P1, P2-A).

349   Nanopore sequencing and assembly (blue) demonstrates better assembly contiguity than read

350   cloud (gold) and short read (green) approaches, but produced a smaller overall assembly with

351   fewer complete drafts at the overall sequence coverage obtained. (a) Relative genus-level

352   abundances are shown for a conventional workflow consisting of bead-beating extraction and

353   short read sequencing, as well as the present workflow consisting of high molecular weight DNA

354   extraction and long read sequencing. (b) For all organisms achieving assembly N50 of at least

355   500 kbp by any approach, genome draft quality and contiguity are shown for long reads, read

356   clouds and short reads. Shapes indicate draft quality. Circularized genomes are indicated by

357   green circles. (c) Complete genome bins with a minimum N50. (d) Complete genome bins below

358   a given read coverage depth. Genome bins with lower read coverage originate from less

359   abundant organisms. (e) Complete genome bins with N50 of > 2 Mbp below a given read

360   coverage depth. (f) High quality genome bins with a minimum N50. (g) High quality genome bins

361    below a given depth of read coverage. (h) High quality genome bins with an N50 exceeding 2

362    Mbp below a given read coverage depth.

363

364    <u>Figure 2</u>

365    Genome assemblies, repeat structure and relative insertion sequence strain abundances of

366    *Prevotella copri* and genome assembly comparisons for other closed genome assemblies. The

367    *Prevotella copri* genome is difficult to assemble beyond insertion sequence sites due to their

368    repetitiveness. For this reason, short read (green) and read cloud (gold) assemblies are highly

369    fragmentary despite very high coverage (>4000x coverage depth). Long reads achieve a closed

370    genome in spite of much lower coverage (318x) (blue). Relative abundances of strains carrying

371    each insertion sequence instance are shown for 0-month and 15-month timepoints (first and

372    second tracks), as well as log-fold change at each site between the two timepoints (third track).

373    b) Finished genomes assembled by the present workflow (blue) are shown with corresponding

374    bins obtained from read cloud (gold) and short read (green) approaches. Read cloud and short

375    read approaches yield more fragmentary approaches, with large genomic regions missing due

376    to incomplete binning.

377   **Supplementary Figure Legends**

378   Supplementary Figure 1

379       Overview of the molecular and informatic workflow steps. a) Extraction consists of

380   enzymatic degradation of bacterial cell walls followed by an initial DNA extraction in phenol-

381   chloroform. This is followed by a proteinase K and RNase A digestion at high temperature and

382   purification with a gravity column. Finally, small fragments are removed by modified SPRI bead

383   size selection. b) After sequencing and basecalling, read sequences are assembled twice with

384   varying genomeSize parameter values. These two assemblies are merged, then circular

385   sequences are identified and trimmed. The consensus sequence is refined by either short-read

386   or long-read polishing.

387

388   Supplementary Figure 2

389       Histogram of total bases versus read length for the three samples sequenced with the

390   current approach.  Read lengths vary between <1kbp to >100kbp, with N50 values between

391   5kbp and 10kbp.

392

393   Supplementary Figure 3

394       Reference alignment dotplots for closed genomes obtained by nanopore long read

395   sequencing and assembly.  Although assemblies share broad structural similarity to available

396   references, there are cases where observed organisms are significantly structurally diverged

397   (e.g. *Dialister*) and in one case bears minimal similarity to the available reference

398   (*Subdoligranulum*).

399

400   Supplementary Figure 4

401       Homopolymer count as a function of length, and homopolymer error in assembled

402   sequence as a function of length in corrected sequence. We found that uncorrected long read

403    assembly demonstrated a 3% error rate with 3-mer homopolymers, assembled too short by an

404    average of 0.5 nucleotides.  This worsens to a 65% error rate on 6-mer homopolymers, which

405    were assembled too short by an average of 1.3 nucleotides.  On average, 63 homopolymers of

406    length 3 or greater were found per kilobase of assembled sequence, of which 4.5 (7.1%) were
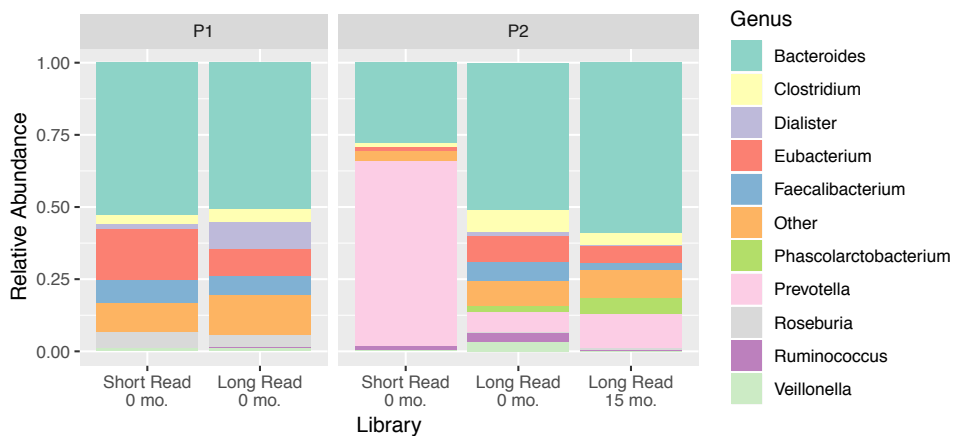
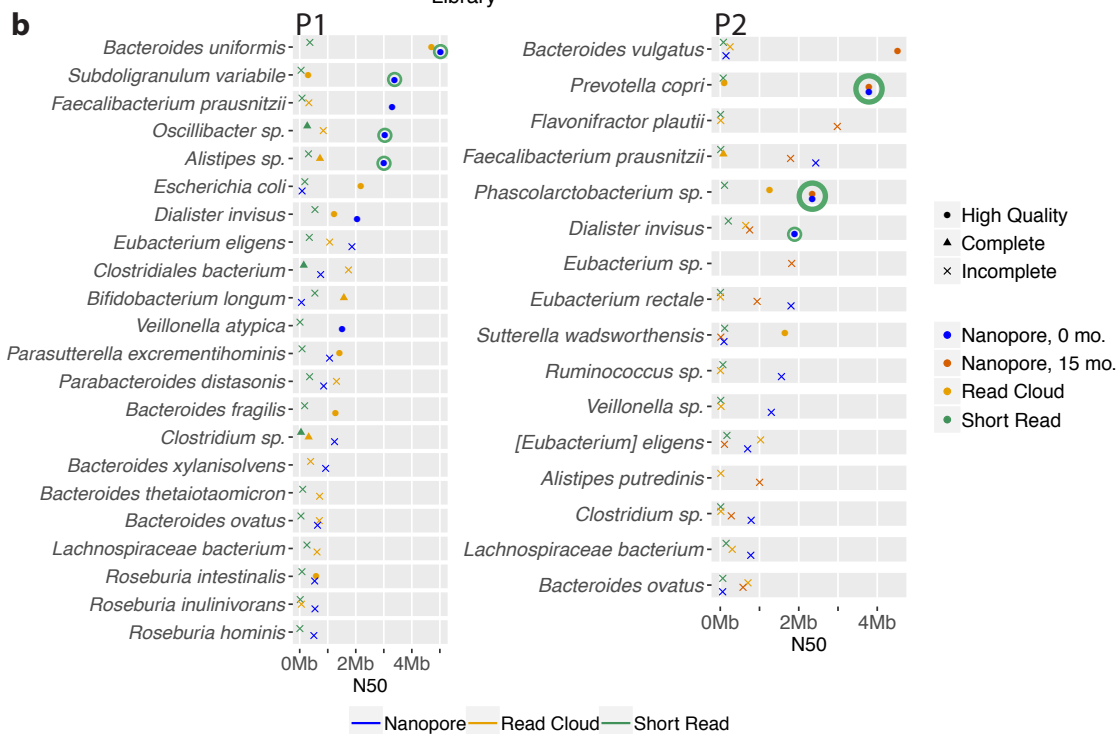407    found to require correction with short reads.

408

409    Supplementary Figure 5

410         Nanopore long read assembly in some cases produces over-circularized genomes.

411    These are sequences that are assembled beyond the wrap-around point, resulting in (a)

412    redundant sequence which are detected and trimmed with the Encircle utility (present study).

413    These sequences can be visualized as (b) corner-cutting off-diagonal alignments within contig

414    self-alignment dotplots, such as that shown for the untrimmed *Subdoligranulum variabile*
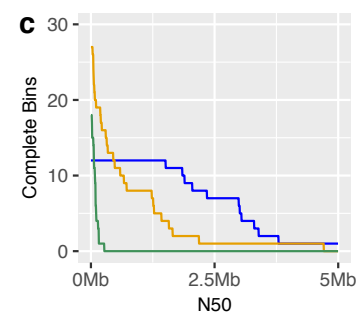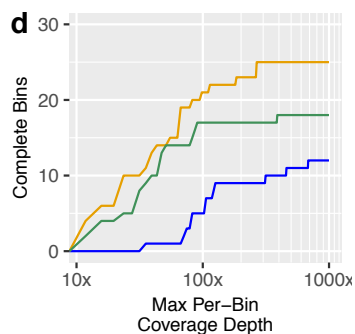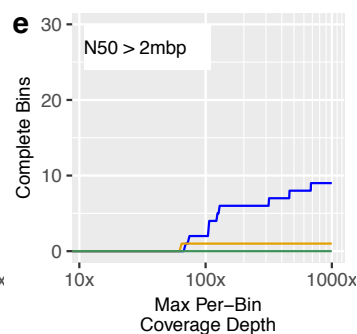
415    assembly.

# Figure 1

# Figure 2



a

P2, 0 Months IS Rel. Abundance

P2, 15 Months IS Rel. Abundance

Log Fold Change

Short Read

Read Cloud

Long Reads

1 mb

2 mb

3 mb

Prevotella copri

b

Oscillibacter sp.

Phascolarctobacterium sp.

Subdoligranulum variabile

Alistipes finegoldii

Bacteroides uniformis

Dialister invisus