1 **TITLE PAGE**

2 **Title**

3 Complete genome assembly of clinical multidrug resistant *Bacteroides fragilis* isolates enables
4 comprehensive identification of antimicrobial resistance genes and plasmids.

5 **Running title**

6 Genomes and plasmids of six clinical MDR *B. fragilis*

7 **Author names**

8 Thomas V. Sydenham[1,2,3]*, Søren Overballe-Petersen[4], Henrik Hasman[4], Hannah Wexler[5], Michael Kemp[1,2],
9 Ulrik S. Justesen[1,2]

10 **Affiliations**

11 [1]Research Unit of Clinical Microbiology, Department of Clinical Research, University of Southern Denmark,
12 Odense, Denmark; [2]Department of Clinical Microbiology, Odense University Hospital, Odense, Denmark;
13 [3]Department of Clinical Microbiology, Lillebaelt Hospital, Vejle, Denmark; [4]Bacteria, Parasites & Fungi,
14 Statens Serum Institut, Copenhagen, Denmark;[5] GLAVA Health Care System and David Geffen School of
15 Medicine at UCLA, Los Angeles, CA, United States.

16 **Corresponding author**

17 Thomas V. Sydenham, Thomas.sydenham@rsyd.dk

18 **ORCIDs**

19 Thomas V. Sydenham: 0000-0003-1058-2449

20 Henrik Hasman: 0000-0001-6314-2709

21 Michael Kemp: 0000-0001-5989-0421

22 Ulrik S. Justesen: 0000-0002-6130-1902

23 **KEYWORDS**

24 *Bacteroides fragilis;* antimicrobial resistance; genome sequencing; plasmid; oxford nanopore; hybrid
25 assembly; insertion sequences

26 **REPOSITORIES**

27 Sequence files (MinION reads de-multiplexed with Deepbinner and basecalled with Albacore in fast5 format
28 and Illumina MiSeq reads in fastq format) and final genome assemblies have been deposited to
29 NCBI/ENA/DDBJ under Bioproject accessions PRJNA525024, PRJNA244942, PRJNA244943, PRJNA244944,
30 PRJNA253771, PRJNA254401, and PRJNA254455

31 **ABSTRACT**

32 *Bacteroides fragilis* constitutes a significant part of the normal human gut microbiota and can also act as an
33 opportunistic pathogen. Antimicrobial resistance and the prevalence of antimicrobial resistance genes are
34 increasing, and prediction of antimicrobial susceptibility based on sequence information could support
35 targeted antimicrobial therapy in a clinical setting. Complete identification of insertion sequence (IS)
36 elements carrying promoter sequences upstream of resistance genes is necessary for prediction of
37 antimicrobial resistance. However, *de novo* assemblies from short reads alone are often fractured due to

1

38  repeat regions and the presence multiple copies of identical IS elements. Identification of plasmids in
39  clinical isolates can aid in the surveillance of the dissemination of antimicrobial resistance and
40  comprehensive sequence databases support microbiome and metagenomic studies. Here we test several
41  short-read, hybrid and long-lead assembly pipelines by assembling the type strain *B. fragilis* CCUG4856T
42  (=ATCC25285=NCTC9343) with Illumina short reads and long reads generated by Oxford Nanopore
43  Technologies (ONT) MinION sequencing. Hybrid assembly with Unicycler, using quality filtered Illumina
44  reads and Filtlong filtered and Canu corrected ONT reads produced the assembly of highest quality. This
45  approach was then applied to six clinical multidrug resistant *B. fragilis* isolates and, with minimal manual
46  finishing of chromosomal assemblies of three isolates, complete, circular assemblies of all isolates were
47  produced. Eleven circular, putative plasmids were identified in the six assemblies of which only three
48  corresponded to a known cultured *Bacteroides* plasmid. Complete IS elements could be identified upstream
49  of antimicrobial resistance genes, however there was not complete correlation between the absence of IS
50  elements and antimicrobial susceptibility. As our knowledge on factors that increase expression of
51  resistance genes in the absence of IS elements is limited, further research is needed prior to implementing
52  antimicrobial resistance prediction for *B. fragilis* from whole genome sequencing.

53  **IMPACT STATEMENT**

54  Bacterial whole genome sequencing is increasingly used in public health, clinical, and research laboratories
55  for typing, identification of virulence factors, phylogenomics, outbreak investigation and identification of
56  antimicrobial resistance genes. In some settings, diagnostic microbiome amplicon sequencing or
57  metagenomic sequencing directly from clinical samples is already implemented and informs treatment
58  decisions. The prospect of prediction of antimicrobial susceptibility based on resistome identification holds
59  promises for shortening time from sample to report and informing treatment decisions. Databases with
60  comprehensive reference sequences of high quality are a necessity for these purposes. *Bacteroides fragilis*
61  is an important part of the human commensal gut microbiota and is also the most commonly isolated
62  anaerobic bacterium from non-faecal clinical samples but few complete genome assemblies are available
63  through public databases. The fragmented assemblies from short read de novo assembly often negate the
64  identification of insertion sequences upstream of antimicrobial resistance gens, which is necessary for
65  prediction of antimicrobial resistance from whole genome sequencing. Here we test multiple assembly
66  pipelines with short read Illumina data and long read data from Oxford Nanopore Technologies MinION
67  sequencing to select an optimal pipeline for complete genome assembly of *B. fragilis*. However, *B. fragilis* is
68  a highly plastic genome with multiple inversive repeat regions, and complete genome assembly of six
69  clinical multidrug resistant isolates still required minor manual finishing for half the isolates. Complete
70  identification of known insertion sequences and resistance genes was possible from the complete genome.
71  In addition, the current catalogue of *Bacteroides* plasmid sequences is augmented by eight new plasmid
72  sequences that do not have corresponding, complete entries in the NCBI database. This work almost
73  doubles the number of publicly available complete, finished chromosomal and plasmid *B. fragilis* sequences
74  paving the way for further studies on antimicrobial resistance prediction and increased quality of
75  microbiome and metagenomic studies.

76  **ABBREVIATIONS**

77  AMR, antimicrobial resistance; WGS, whole genome sequencing; IS, insertion sequence; ONT, Oxford
78  Nanopore Technologies;

79  **DATA SUMMARY**

80  1.  Sequence read files (Oxford Nanopore (ONT) fast5 files and Illumina fastq files) as well as the final

81      genome assemblies have been deposited to NCBI/ENA/DDBJ under Bioproject accessions

82      PRJNA525024, PRJNA244942, PRJNA244943, PRJNA244944, PRJNA253771, PRJNA254401, and

83      PRJNA254455.

84   2.  Fastq format of demultiplexed ONT reads trimmed of adapters and barcode sequences are available at

85      doi.org/10.5281/zenodo.2677927

86   3.  Genome assemblies from the assembly pipeline validation are available at doi:

87      doi.org/10.5281/zenodo.2648546.

88   4.  Genome assemblies corresponding to each stage of the process of the assembly are available at

89      doi.org/10.5281/zenodo.2661704.

90   5.  Full commands and scripts used are available from GitHub: https://github.com/thsyd/bfassembly

91      as well as a static version at doi.org/10.5281/zenodo.2683511

## INTRODUCTION

93   *Bacteroides fragilis* is a Gram-negative anaerobic bacterium that is commensal to the human gut but can
94   act as an opportunistic pathogen; it is the most commonly isolated anaerobic bacteria from non-faecal
95   clinical samples (1). Antimicrobial resistance rates are increasing for *B. fragilis*, especially for carbapenems
96   and metronidazole, two widely used antimicrobials for treatment of severe infections and anaerobe
97   bacteria (2,3). Antimicrobial susceptibility testing of anaerobes using agar dilution or gradient strip methods
98   can be costly and labour intensive and despite efforts to validate disk diffusion as a less expensive option,
99   turn-around time will still be least 18 hours and validation for individual species will be required (4).

100  Antimicrobial resistance prediction from bacterial whole genome sequences, from cultured isolates as well
101  as metagenomes, could be implemented in clinical microbiology in the near future, with the potential for
102  improved sample-to-report turnover time and possibly eliminating the need for phenotypical testing for
103  individual species (5–8). For a few species, prediction of antimicrobial resistance from WGS has been
104  validated, but for the majority of clinical relevant species challenges still remain (6,9,10).

105  Based on DNA-DNA hybridisation studies, *B. fragilis* can be divided into two DNA homology groups (division
106  I and II), whose ribosomal contents are so different that the two divisions can be distinguished by mass
107  spectrometry routinely used to identify isolates in clinical laboratories (11). *B. fragilis* division I carry the
108  chromosomal cephalosporinase gene *cepA* whilst *B. fragilis* division II harbour the chromosomal metallo-β-
109  lactamase gene *cfiA* (also known as *ccrA*) (12,13). The *cfiA* gene can confer resistance to carbapenems, a
110  class of antimicrobials usually reserved for patients with severe sepsis or infections with multidrug-resistant
111  bacteria. But expression levels are partly controlled by insertion sequence (IS) elements carrying promotor
112  sequences inserted upstream of the gene and only 30-50% of clinical isolates that harbour *cfiA* display
113  phenotypically reduced susceptibility to carbapenems (3). The same pattern of expression control can be
114  observed for genes associated with resistance to metronidazole (*nim* genes) and clindamycin (*erm* genes)
115  (1).

116  In 2014 we observed that identification of IS elements upstream of known antimicrobial resistance genes in
117  *B. fragilis* was hampered in short read *de novo* assemblies even though the genes could be identified (14).
118  This occurred because contigs were often terminated close to the start of the resistance genes, presumably

119  due to the proliferation of multiple copies of the same IS elements throughout the *B. fragilis* genomes.
120  Genome assemblies from short read sequencing technologies alone most often result in fragmented
121  assemblies because of repetitive regions and genome elements with multiple occurrences in the
122  chromosomes and plasmids (15,16). Therefore, we could not predict antimicrobial resistance (AMR)
123  phenotypes in *B. fragilis* using only short reads for WGS since IS element identification is a prerequisite for
124  correct genotype-phenotype associations. Long read sequencing technologies are increasingly being
125  utilised to increase the contiguity of bacterial genome assemblies and often result in complete, closed
126  chromosomes and plasmids (17–20). This provides possibilities for comprehensive identification of IS
127  elements, insights into genome structures and characterisation of other mobilisable elements and
128  associated genes. Complete identification and characterisation of plasmids in sequenced isolates would
129  allow for improved analysis of the plasmid-mediated spread of antimicrobial resistance.

130  Bioinformatic analysis of WGS data depends heavily on high-quality reference databases. Anaerobes make
131  up most of the bacterial human commensal microbiota but are most likely underrepresented in public
132  databases of whole genomes from cultured isolates. The NCBI Genome database (accessed 31-03-2019)
133  contains genome sequences of 191,411 bacteria of which 13,483 are marked as complete assemblies. Only
134  seven of these are *Bacteroides fragilis* (21–27). In comparison there are 776 assemblies of *E. coli* marked as
135  complete and 398 of *S. aureus*. Improving the representation of complete assemblies of *B. fragilis* in the
136  public genome databases will support the development of antimicrobial resistance prediction from WGS as
137  well as microbiome and metagenomic analysis projects.

138  The aims of this study were to select an optimal assembly software pipeline for complete, circular assembly
139  of *Bacteroides fragilis* and demonstrate the utility of complete assembly for both plasmid identification and
140  comprehensive detection of genes and IS elements associated with antimicrobial resistance. We assembled
141  the *B. fragilis* CCUG4856T (= ATCC25285 = NCTC9343) reference strain utilising long reads generated with
142  the MinION sequencer from Oxford Nanopore Technology (ONT) and high-quality Illumina short reads and
143  selected the best assembly pipeline by comparing assemblies to the Sanger sequenced reference NCTC9343
144  (RefSeq accesion GCF_000025985.1). The best assembly pipeline was then applied to six clinical multi-drug
145  resistant *B. fragilis* isolates from our 2014 study (14).

146  **METHODS**

147  **Culture conditions and DNA extraction**

148  *Bacteroides fragilis* CCUG4856T and the six strains described in our previous study were included (14,21).
149  Strains were stored at -80° in beef extract broth with 10% glycerol (SSI Diagnostica) and cultured on solid
150  chocolate agar with added vitamin K and cysteine (SSI Diagnostica) for 48 hrs in an anaerobic atmosphere
151  at 35 °C. Ten µl of culture was transferred to 14 ml saccharose serum broth (SSI Diagnostica) and incubated
152  for 18 hrs under the same conditions. DNA was then extracted using the Genomic-Tip G/500 kit (Qiagen)
153  following the manufacturers protocol for Gram negative bacteria and eluted into 5 mM Tris pH 7.5 0.5 mM
154  EDTA buffer. Quality control was performed by measuring fragment length on a TapeStation 2500
155  (Genomic DNA ScreenTape, Agilent), purity on the NanoDrop (ThermoFisher Scientific) and concentration
156  on the Qubit (dsDNA BR kit; Invitrogen). The eluted DNA was then stored at -20 °C.

157  **Illumina library preparation, sequencing and quality control**

158  The strains had previously been sequenced and assembled using Illumina short reads for our previous study
159  (14), but to minimise biological disparities we opted to re-sequence with Illumina using the same DNA
160  extraction prepared for long read sequencing. Paired-end libraries were generated using the Nextera XT
161  DNA sample preparation kit (Illumina) according to the manufacturer's protocol. DNA was sequenced on a
162  MiSeq sequencer (Illumina) with 150 bp reads for a theoretical read depth of 100x. Read quality metrics
163  were evaluated using FastQC (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) and fastp

164  v0.19.6 (28). Filterbytile from the BBmap package (http://sourceforge.net/projects/bbmap/) was used for
165  removing low-quality reads based on positional information on the sequencing flowcell and TrimGalore
166  (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/), with settings --qual 20 and --length
167  126, provided additional adapter and quality trimming. FastQ files were then randomly down-sampled to <
168  100x crude read depth using an estimated genome size of 5.3 Mb, as higher read depths tend to reduce
169  assembly quality (29).

**Nanopore library preparation and MinION sequencing**

171  Sequencing libraries were prepared using the Rapid Barcoding kit (SQK-RPB004; Oxford Nanopore
172  Technologies) following the manufacturers protocol (version RPB_9059_v1_revC_08Mar2018) with SPRI
173  bead clean up (AMPure XT beads; Beckman Coulter) as described. Sequencing was performed as multiplex
174  runs on a MinION connected to a Windows PC with MinKnow v1.15.1 using FLO-MIN106 R9.4 flowcells.
175  Raw fast5 files were transferred to the Computerome high performance cluster
176  (https://www.computerome.dk/) for analysis. Four sequencing runs were performed, as the first two runs
177  did not provide enough data for complete assembly of all isolates (see results section).

**Fast5 demultiplexing, base-calling, quality control and filtering**

179  The raw fast5 files were demultiplexed with Deepbinner v0.2.0 and base-called using Albacore v2.3.3,
180  retaining only those barcodes Deepbinner and Albacore agreed upon for minimal barcode misclassification
181  (30). Porechop v0.2.4 (https://github.com/rrwick/Porechop) with the --*discard_middle* option was used for
182  adapter and barcode trimming and read statistics were collected using NanoPlot (31). Filtlong v0.2.0
183  (https://github.com/rrwick/Filtlong) was used to filter the long reads by either removing the worst 10% or
184  by retaining 500Mbs in total, which ever option resulted in fewer reads.

**Assembly validation**

186  To select and validate the optimum assembly pipeline *Bacteroides fragilis* CCUG4856T was assembled using
187  a variety of well-known assemblers and polishing tools (Table 1). Each assembler was run with the Filtlong
188  filtered reads as input or the filtered reads corrected with Canu 1.8 (with standard settings,
189  corMinCoverage=0, or coroutCoverage=999). Canu was also tested with the unfiltered reads as input.
190  Hybrid assemblers used the filtered long reads and the filtered, trimmed and down-sampled Illumina reads.
191  Unicycler includes polishing with Racon and Pilon. For assemblers other than Unicycler, Racon polishing
192  with ONT reads was run for one or two rounds and Pilon was run until no changes were made or for a
193  maximum of six rounds. Racon polishing with Illumina reads was run for one round.

194  The original Sanger sequenced *Bacteroides fragilis* NCTC9343 (=CCUG4856T) (21) downloaded from NCBI
195  RefSeq (accession GCF_000025985.1) was used as reference sequence for the assembly comparisons and
196  Quast v5.0.2 was used for assembly summary statistics, indel count, and K-mer-based completion (32).
197  BUSCO v3.0.2b with the bacteroidetes_odb9 dataset, CheckM v1.0.12, and Prokka v1.13.3 were used to
198  assess gene content (33–35). Average nucleotide identity was calculated using
199  https://github.com/chjp/ANI/blob/master/ANI.pl and ALE v0.9, which uses a likelihood based approach to
200  assess the quality of different assemblies, was also used to score the assemblies (36,37). Ranking of
201  assemblies was based on number of contigs, number of circular contigs, closeness to total length compared
202  to the reference genome, number of local misassembles, number of mismatches per 100 kb, number of
203  indels per 100kb, average nucleotide identity (ANI), CheckM and BUSCO scores, and the total ALE score (a
204  higher score is better). Please see https://github.com/thsyd/bfassembly for full bioinformatics methods.

**Genome assembly of MDR *B. fragilis* isolates**

206  The assembly strategy deemed to produce the highest quality genome for CCUG4856T was chosen for
207  initial assembly of the six MDR *B. fragilis* isolates. Manual finishing of incomplete assemblies was

208  performed using Bandage for visualisation of assembly graphs and BLASTn searches (38). Minimap2 and
209  BWA MEM were used to map reads to the assemblies for coverage graphs (39,40). Long read assembly with
210  Flye was compared to the Unicycler assembly and used to guide and validate the manual finishing results.
211  Circlator's *fixstart* task was used to fix the start position of the manually finished genomes to be at the
212  *dnaA* gene (41).

213  The assembled genomes were submitted to NCBI GenBank and annotated with PGAP (42). ABRicate v0.8.10
214  (https://github.com/tseemann/ABRicate) (with options --minid 40 --mincov 25) was used to screen for
215  antimicrobial resistance genes with the ResFinder (database date 19-08-2018), NCBI Bacterial Antimicrobial
216  Resistance Reference Gene Database (database date 19-09-2018), and CARD (v2.0.3) databases,
217  supplemented with nucleotide sequences for the multidrug efflux-pump genes *bexA* (GenBank:
218  AB067769.1:3564..4895) and *bexB* (GenBank: AY375536.1:4599..5963) (43,44). IS elements were identified
219  using ABRicate with data from the IS-finder database (http://www-is.biotoul.fr/, Update: 2018-07-25) (45).

**Identification of plasmids and mobile genetic elements**

221  The PLSDB web server (https://ccb-microbe.cs.uni-saarland.de/plsdb/) (data v. 2019_03_05) contains
222  bacterial plasmid sequences retrieved from the NCBI and was used for screening  and identifying putative
223  plasmids sequences (46). Only hits to accessions from cultured organisms were included. Putative plasmids
224  not identified using PLSDB, were evaluated by the read depth relative to the chromosome (higher relative
225  read depth indicates plasmid sequence) and Pfam families covering known plasmid replication domains
226  from Table 1 in reference (47) were downloaded from the Pfam database (Pfam 32.0,
227  https://pfam.xfam.org/) and used for screening putative plasmids with ABRicate.

**RESULTS**

**Sequencing data quality**

230  For Illumina data, a median of 3,465,082 reads (interquartile range [IQR]: 3,177,493-5,001,077) were
231  generated for each isolate (Supplementary Table S1)). After filtering, adapter-removal and down sampling a
232  median of 449,022,741 bases (IQR: 433,517,549-530,257,210) were available per isolate with 87-96% Q30
233  bases corresponding to calculated read depths of 75-103%. The %GC content of the reads for each isolate
234  (median 42.9%, range: 42.6-43.3%) were very consistent and within the expected range for the *Bacteroides*
235  genus (40-48%) (48) .

236  Isolates were sequenced in runs multiplexed with other isolates not included in this study. Based on initial
237  test assemblies using Unicycler without filtering or Canu correction (not shown) it was concluded that data
238  from the first ONT sequencing runs were to be supplemented by additional runs to increase the chance of
239  complete assembly of all isolates. Concatenating reads from runs, a median of 75,598 reads [IQR: 50,210-
240  112,065] with a median length of 2,938-4,393 bases were generated for each isolate (Supplementary Table
241  S1). Filtering with Filtlong and correction with Canu resulted in a median of 8,515 reads (IQR: 6,226-10,370)
242  with median lengths of 6,181-38,588 for each isolate as input for the assemblies.

**Selecting the optimal assembly pipeline**

244  141 assemblies of *B. fragilis* CCUG4856T were generated using the various assemblers and polishing steps
245  (Supplementary Table S2). Compared to the reference genome, Unicycler assemblies were of the highest
246  quality (Table 2**Error! Reference source not found.**). Unicycler, with any of the read input options,
247  produced two circular contigs of the expected lengths, and the differences between the various Unicycler
248  assemblies were minimal (Table 3**Error! Reference source not found.**). Assemblies with Canu corrected
249  reads showed slightly higher genome fractions and average nucleotide identities to the reference and
250  fewer mismatches and indels, when compared to Unicycler alone. Unicycler assemblies corrected with
251  Racon using Illumina reads worsened slightly overall with 0.04-0.19 more indels and 0.14-0.25 more

252  mismatches per 100 kbp. Based on this initial evaluation, the assembly pipeline using Canu corrected reads
253  with default options was chosen (Assembly "OF.CS" in Table 3). This would reduce the number of long
254  reads, compared to Canu correction with corMinCoverage=0, or coroutCoverage=999, and thereby lead to
255  a faster run-time for Unicycler.

256  The hybrid Unicycler assembly of CCUG4856T with standard Canu corrected ONT reads consists of two
257  circular contigs of 5,205,133 and 36,560 bp in length. The plasmid is the same length as plasmid pBF9343
258  from the reference assembly GCF_000025985.1 and the chromosome is seven bases shorter. Alignments of
259  the Sanger sequenced assembly GCF_000025985.1 with the hybrid Unicycler assembly show an 88,045 bp
260  inversion in the hybrid assembly compared to the Sanger assembly (Figure 1). This inversion is present in all
261  the best assemblies, including assemblies derived from solely ONT sequences or Illumina sequences
262  (Supplementary Figure S1) as well as two additional assemblies of NCTC9343/ATCC25285 from PacBio and
263  Illumina sequences downloaded from NCBI RefSeq (Supplementary Figure S2).

264  **Complete assembly of six multidrug resistant isolates**

265  Unicycler, using filtered and trimmed Illumina reads and the Filtlong filtered and Canu corrected ONT reads
266  from the first sequencing runs, generated complete, continuous, circular assemblies for two of the six

7

267     isolates (BFO18 and BFO67) (



| Strain | Best Spades graph | Hybrid assembly after first nanopore sequencing run. No Canu correction | Hybrid assembly after first nanopore sequencing run with Canu correction | Hybrid assembly with data from supplementary runs with Canu correction | After manual finishing | Final assemblies contig statistics | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | No. | Length | Relative read depth |
| CCUG 4856T | | | | | | 1 2 | 5,205,133 36,560 | 1.00x 7.42x |
| BFO17 | | | | | | 1 2 3 | 5,474,541 85,671 5,594 | 1.00x 1.85x 23.03x |
| BFO18 | | | | | | 1 2 3 4 | 5,302,644 7,221 4,137 2,782 | 1.00x 25.98x 50.80x 59.91x |
| S01 | | | | | | 1 2 3 4 | 5,325,251 78,085 8,331 5,595 | 1.00x 2.29x 20.99x 22.67x |
| BFO42 | | | | | | 1 2 3 | 5,141,257 8,306 5,594 | 1.00x 40.06x 40.15x |
| BFO67 | | | | | | 1 2 | 5,478,614 6,129 | 1.00x 94.15x |
| BFO85 | | | | | | 1 | 5,504,076 | 1.00x |

268

269 Figure 2). For the assemblies that were not complete with sequencing data from the first MinION runs,
270 increasing the amount of ONT data resulted in fewer contigs overall, except for BF067, where the additional
271 data from the second sequencing run led to a fragmented assembly and manual finishing was necessary.
272 Performing assembly of isolate S01 without Canu correction of the ONT reads from the first sequencing
273 resulted in a closed chromosome and performing Canu correction of reads resulted in a fragmentation of
274 the chromosome. This was ameliorated by including more ONT data. By manual finishing using read
275 mapping and additional assembly with Flye, the remaining three assemblies were circularised.

8

276    Chromosomes varied in length from 5,141,257 – 5,504,076 bp. Alignment of ONT and Illumina reads to the
277    chromosome assemblies showed even coverage for both sequencing technologies (Supplementary Figure
278    S3). For BFO85 a >100% relative read depth increase was observed at approximately 25kb-38kb. This could
279    represent a 12 kb repeat region that was not resolved in the assembly. Seven (47%) of the 15 PGAP
280    annotated CDS' in the 13kb region were annotated as hypothetical proteins. None of the annotated CDS'
281    represented mobilisable proteins.

**Eleven putative plasmid sequences were identified**

283    A total of 11 putative circular plasmids were identified in the six *B. fragilis* isolates (Table 4). Zero to three
284    putative plasmids were identified per isolate with lengths varying from 2,782 to 85,671bp.

285    The PLSDB database contains NCBI RefSeq plasmid sequences marked as complete. Three of the 11
286    putative plasmid sequences were found to match (ID > 98%) a sequence in PLSDB (Table 4). These three all
287    matched the cryptic plasmid pBFP35 (49). The NCBI Nucleotide database was queried using BLASTn with
288    the remaining unidentified putative plasmid sequences (50). BFO18 putative plasmid sequence pBFO18_1
289    (7,221 bp) resembles plasmid pIP421, a 7.2kb plasmid with metronidazole resistance gene *nimD* and
290    IS*1169.* Partial sequences in NCBI GenBank spanning the *nimD* gene, IS element and *RepA* (GenBank
291    Y10480.1 and X86702.1) showed 99 %ID to their alignment to pBFO18_1 (not shown) (51,52). Strain S01
292    putative plasmid sequence pBFS01_2 (8,331 bp) showed 99.87 %ID to the 1486bp partial sequence of *B.*
293    *fragilis* plasmid pBF388c (GenBank AM042593.1), a 8.3kb conjugative plasmid harbouring *nimE* and IS*Bf6*
294    (53).

295    None of the three putative plasmid sequences of strain BFO18 could be identified using the PLSDB but
296    querying the NCBI nucleotide database using BLASTn revealed hits for all three. The hits corresponded to
297    circularised sequences (%ID: 99.56-99.96, %COV: 100) assembled from mobilome metagenomic sequencing
298    of the uncultured caecum content from a rat trapped at Bispebjerg Hospital in Copenhagen, Denmark (two
299    hours' drive from Odense University Hospital where BFO18 was isolated from a patient's blood culture)
300    (Supplementary Table S3) (47,54). BLASTn searches of the remaining unidentified putative plasmids from
301    the other strains did not reveal complete hits.

302    Using ABRicate with the plasmid replication domains collected from the Pfam database, all putative
303    plasmids, except pBF017_1 and pBFS01_1, were found to have recognised replicon domains (Table 4). DNA
304    fragments of sizes matching pBFO17_1 and pBFS01_1 were detected by PFGE of S1 endonuclease
305    restriction enzyme treated plasmid DNA extracts (Supplementary Figure S4) and the circular structures of
306    the two sequences lacking a predicted replication domain, were confirmed manually by visually inspecting
307    BLASTn mapping of ONT sequences longer than 10 kbp to the assembled plasmid sequences with CLC
308    Genomics Workbench 10 (Qiagen). Eleven and 22 ONT reads spanned the complete lengths of pBFO17_1
309    and pBFS01_1 respectively and contained no other elements. pBFO17_1 and pBFS01_1 demonstrate a
310    degree of similarity of close to 100%, except for an approximate total of 7,500 bp transposase and
311    prophage sequences in pBF017_1 (Figure 3). No alignment to chromosomal sequences of any of the
312    included *B. fragilis* isolates was observed using progressiveMauve (not shown) (55).

313    The GC content of pBFO17_1 and pBFS01_1 are 36.78% and 36.04% respectively. These lie within the range
314    for the *Bacteroides* genus but differ from the expected value for *B. fragilis* (43%), which could indicate that
315    the putative plasmids do not originate from *B. fragilis* (56). After supplementing the PGAP annotations with
316    RAST annotation (57), 63% (pBFO17_1) and 59% (pBFSO1_1) of CDS' remained annotated as hypothetical
317    or as domain of unknown function. Of the annotated CDS' the majority were associated with mobilisable
318    features, plasmids and phages such as *parA* and *parB*, DNA partitioning proteins, conjugative transposon
319    proteins, transposases, DNA binding motif domain containing proteins, and reverse transcriptase protein.
320    The results above support the assembly data suggesting these two sequences are in fact plasmids.

**Detection of antimicrobial resistance genes and Insertion sequence elements**

We used ABRicate to screen assemblies for AMR genes (ResFinder, NCBI and CARD databases supplemented with sequences for *bexA* and *bexB*) and IS elements (IS-finder database); several AMR genes, possible homologs to known AMR genes and IS elements adjunct to the AMR genes were detected (Table 5). Of note, isolate BFO17 contains two homologs of the metronidazole resistance gene *nimJ* (with a 100% consensus) and two isolates, S01 and BFO85, harbour two homologs of the tetracycline resistance gene *tetQ*. Homologs to *bexA* and *bexB* were identified with 73.53-99.12 %ID and were all confirmed with BLASTx searches against the NCBI nr database, as was done in our previous study (14). Partial hits for *ugd* was observed for several isolates, but with low %ID and %COV, and possible represent identification of conserved domains, but not *ugd* homologs. Increased expression of the *cfiA* metallo-beta-lactamase gene, *nim*-family 5-nitroimidazole genes and *erm* genes is partly regulated through IS elements containing promoter sequences. Full length IS elements could be identified upstream of 11 (79%) of 14 *cfiA*, *nim* and *erm* genes and upstream of two of three *CfxA4* genes and the *OXA-347* gene identified in BFO42. The described *Bacteroides fragilis* promotors TAnnTTTG (-7) and TG or TTG or TGTG (-33) (58) were searched for manually, but could not be identified upstream of the two *cfiA* genes in isolates BFO67 and BFO85 or the *ermB* gene in BFO85 for which no IS elements could be detected upstream (not shown).

**Correlation between identified genes and IS elements and phenotypical resistance**

As in our previous study, the *cfiA* gene was identified in the five meropenem resistant isolates (Table 5). All the *cfiA* genes were found on the chromosomal sequences. Complete IS elements were identified upstream of the *cfiA* genes in BFO17, BFO18 and S01, but not in BFO67 or BFO85. MICs for meropenem and imipenem were lower for these two isolates. *Nim* genes (-A, -D, -E and -J) could be found in the four metronidazole resistant isolates, all with complete IS elements upstream. Three of the *nim* genes were found on putative plasmids of the respective isolates. The four clindamycin-resistant isolates all carried *erm*-genes with upstream IS elements. A transposase was inserted in the *ermF-gene* in isolate BFO18, splitting it in two and the same isolate demonstrated a lower clindamycin MIC (6 mg/L) than the other three clindamycin resistant isolates.

**DISCUSSION**

**Hybrid genome assembly produces high quality *B. fragilis* genomes**

The primary aim of this study was to select and validate an assembly method to reliably complete chromosome and plasmid assembly of *B. fragilis* genomes. From 141 assembly variations, a hybrid approach using Filtlong filtered and Canu corrected ONT reads with quality filtered Illumina reads as input to Unicycler produced a complete, closed assembly of *B. fragilis* CCUG4856T with high similarity to the reference assembly of the original Sanger sequenced reference assembly. An 88kb inversion was observed when comparing the two assemblies. Cerdeño-Tárraga and colleagues observed difficulties in resolving certain regions of the Sanger sequenced assembly of NCTC9343 due to invertible regions with flanking inverted repeat sequences (21). The observed inversion in the hybrid Unicycler assembly, could be due to a) a superior assembly where the longer ONT reads have overcome the shortcomings of the shorter Sanger sequences, b) an incorrect assembly by Unicycler, c) a biological difference that has occurred over time between the strain stored at NCTC and CCUG, or d) a biological difference that occurred during the culturing of the strain, with dominance of a clone with the inversion, prior to DNA extraction as part of this study. The observations that the inversion is also present in all the best assemblies from this study and assemblies from two other research institutions support the conclusions that the current hybrid Unicycler assembly represents the true orientation of the 88kb sequence.

**Complete genome assembly of three of the six multidrug resistant isolates required manual finishing**

10

365  The assemblies of BFO18, S01, and BFO42 were completed by Unicycler without manual intervention, but
366  the chromosomes of BFO17, BFO67, and BFO85 could only be closed by performing manual steps. The
367  manual finishing steps are time consuming, difficult to replicate and are easily biased. In order to be
368  implemented in routine clinical laboratories, large scale, automated, complete assembly of prokaryote
369  genomes require robust methods with minimal human interaction. Genome assembly using another long-
370  read assembler, Flye, supported the results of the manual finishing for two of three isolates. Flye is better
371  at resolving repeats than miniasm, the long read assembler included in the Unicycler pipeline (59). One
372  option could be to include the long-read assembly from Flye, in place of that of miniasm, to guide bridge
373  building for the higher quality Illumina-only contigs produced in the first steps of Unicycler. To resolve
374  repeats it is often necessary to have long reads that span the repeat. In prokaryotes repeats over 10kb are
375  not unusual and they are often spanned by the ONT reads generated, even by novice researchers. But
376  repeat regions of up to 120kb and duplications of 200kb have been described in some prokaryotes
377  (17,18,60). ONT sequencing runs will routinely result in many reads that span the majority of repeats, but
378  to obtain ONT reads that span specific 120-200kb repeats in a genome of interest still requires skill and a
379  certain amount of luck. Protocols for ONT sequencing have been described that result in read lengths of
380  over 2 Mb, but this requires skilled and experienced researchers and lab technicians and demands high
381  amounts of very high quality input DNA and essentially sequencing of only one isolate per MinION flowcell
382  (61).

383  ONT read depth did not serve as an indicator of whether the Unicycler assemblies would result in closed
384  chromosomal contigs in this study. Final ONT read depth, prior to Filtlong filtering and Canu correction,
385  ranged from 23-371x, but a high read depth alone, was not an indicator of closed contigs. The three
386  assemblies BF017, BFO67, and BFO85 required manual finishing to complete the assemblies and had ONT
387  raw read depths of 99-137x. After Filtlong filtering and Canu correction the median read lengths were
388  21,932-29,893b and read length N50 was 25,765-34,815b for the three isolates (Supplementary Table S1).
389  Canu correction improved the Unicycler assembly of *B. fragilis* CCUG4856T by nearly all parameters. But
390  whilst Canu correction of the data from the first sequencing run resulted in the complete assembly of
391  BFO67, the assembly of S01 worsened slightly. Increasing the amount of ONT data for BFO67 fragmented
392  the complete chromosome. However, increasing the ONT read depth did decrease the number of contigs
393  per isolate in our study overall.

394  Defining an optimal approach for complete prokaryote genome assembly is a continuous process, as
395  sequencing technologies and assembly software develop and mature. Ring and colleagues found that Canu
396  correction prior to Unicycler hybrid assembly was superior to other hybrid assembly or long read assembly
397  approaches for assembly of *Bordetella pertussis* genomes that contain long duplicated regions (18).
398  Unicycler also performs well in other studies comparing assessing genome assemblers for bacterial genome
399  and plasmid assembly (19). De Maio and colleagues recently published a preprint comparing hybrid
400  assembly strategies for 20 *Enterobacteriaceae* isolates (20). In their dataset, simply randomly subsampling
401  ONT reads to an approximate read depth of 100x was slightly superior to applying Canu correction or
402  Filtlong filtering prior to Unicycler assembly. For 85% of isolates the expected number of circular contigs
403  were all assembled. For only one additional isolate Canu correction or Filtlong filtering resulted in the
404  assembly of the expected number of circular contigs. Manual steps, including down sampling ONT reads or
405  removing the Canu correction are options to consider, if chromosomes are not complete and circularised
406  after initial Unicycler assembly, providing ONT read depth of 100x is available.

407  We chose to benchmark a selection of widely used genome assemblers for short read, long-read and hybrid
408  bacterial genome assembly as well as polishing tools for long read assemblies, but many other options have
409  been published. Most assemblers and polishing tools were run using default parameters, and it is possible
410  that further optimisation of settings for the individual software packages might have improved assemblies
411  further than was demonstrated here. As sequencing technologies and assembly software continues to

11

412 improve, continued validation of pipelines is advisable. Software such as poreTally provides user friendly
413 options for benchmarking genome assembly pipelines prior to implementation (62).

414 ***Bacteroides* plasmids are not well represented in public databases**

415 A secondary aim of this study was to identify plasmids in the hybrid assemblies. Automated tools have been
416 developed and validated for identification of plasmids from genome assemblies or read data, but they are
417 dependant of collated databases of known plasmid sequences. As such, tools such as PlasmidFinder or
418 mlplasmids can be applied for plasmid identification for *Enterobacteriaceae* or *Enterococcus faecium*, but *B.*
419 *fragilis* is not supported at the time of writing (63,64). Therefore, we evaluated putative plasmid sequences
420 by sequence identity and length comparison using the PLSDB webpage, identifying plasmid replication
421 domains, and using circularisation and relative coverage as indicators that a sequence represents a plasmid
422 in a given isolate.

423 Only four of the twelve plasmid sequences from the seven isolates could be identified using the PLSDB and
424 three of these were the same plasmid, pBFP35. Two other putative plasmids, pBFO18_1 and pBFS01_2
425 were likely plasmids pBF388c and pIP421 based on the partial sequences from these plasmids and plasmid
426 length. This still leaves half of the circularised, putative plasmids unidentified. The two longer putative
427 plasmids, pBFO17_1 and pBFS01_1, displayed a high degree of similarity, a GC% out of the normal range for
428 *B. fragilis*, and a relative read depth of double the reads compared to the chromosome. Most annotated
429 CDS' were associated with mobilisable elements, but no known plasmid replication domains could be
430 identified. From the sequencing data alone, we cannot conclude that they represent true plasmids,
431 however the findings above and manual inspection of long read mapping support that inference.

432 There are only 14 complete plasmid sequences from cultured *Bacteroides* isolates in the PLSDB
433 v2019_03_05, which is based on the NCBI RefSeq database. Many other *Bacteroides* plasmids have been
434 partially described, and some are represented by partial sequences or marked as contig level in the NCBI
435 nucleotide database (65–68). Metagenomic sequencing and genome assembly projects are expanding the
436 public sequence databases and screening the NCBI nucleotide database, sequences with a high degree of
437 similarity to the putative plasmid sequences from one patient isolate (BFO18) could be found. These
438 originated from a rat caecum metagenomic plasmid sequencing project from Copenhagen, a few hours'
439 drive from Odense University Hospital. To understand and perform surveillance of the dissemination of
440 plasmids there is a need for increased submissions of high quality, annotated and phenotypically validated
441 sequences of bacterial isolates including plasmids. This study adds significantly to the number of complete
442 plasmid sequences associated with *Bacteroides.*

443 **Complete assembly allows comprehensive identification of resistance determinants in *B. fragilis***

444 We also intended to comprehensively identify resistance genes and IS elements in the hybrid genome
445 assemblies. Using ABRicate with several resistance gene databases and IS-element nucleotide sequences,
446 the findings of our previous study were confirmed and enhanced. Assemblies from Illumina sequencing
447 alone would only allow partial IS element identification (14). Now, with the complete assemblies,
448 comprehensive identification of known IS elements upstream of the relevant resistance genes could be
449 completed. In our first study we used ResFinder with the available database at that time. Now, by including
450 several databases, and lowering the %ID threshold, the number of genes identified increased. Additionally,
451 for as a result of the complete genome assembly of BFO17, we could now identify two copies of *nimJ*, while
452 only one copy was identified in the short read draft assembly of the same isolate in the previous study.
453 Husain and colleagues identified the presence of three copies of *nimJ* in strain HMW615, when describing
454 the *nimJ* gene (69). We confirmed this finding by running ABRicate on the HMW615 assembly as done with
455 the isolates of this study (not shown). Interestingly, RAST annotates a third *nim* gene (nucleotide positions
456 1,359,590..1,360,093) in the Unicycler hybrid assembly of BFO17, and the PGAP annotation includes an
457 additional annotation of a pyridoxamine 5'-phosphate oxidase family gene (nucleotide positions

458 940,032..940,505), the family that includes the *nim*-genes. It is possible that one or more novel homologs
459 of the *nim* are present in BFO17.

460 IS elements could be identified upstream of most relevant resistance genes. However, in three cases no IS
461 element was present upstream of a resistance gene, even though the isolates displayed phenotypical
462 resistance associated with increased expression of the specific gene. Known *B. fragilis* promoter sequences
463 could not be identified upstream of the genes "missing" upstream IS elements, however *B. fragilis*
464 promotors are still not completely described, so it is possible there are other unknown variants.

465 By selecting an optimal genome assembly strategy for *B. fragilis*, supplemented with minimal manual
466 finishing efforts, and applying this to six multidrug resistant isolates, the number of complete *B. fragilis*
467 genomes and plasmids in the public databases has now almost doubled. The future aim of performing
468 antimicrobial resistance prediction based solely on WGS information for *B. fragilis* demands near-complete
469 genomes for identification of IS elements upstream of resistance genes. However, we must caution that the
470 absence of an IS element upstream of *cfiA* does not always correlate to susceptibility to carbapenems.
471 Future studies are needed to address this, and utilising complete genome assembly for genome wide
472 association studies is one approach that could be pursued. Technologies that provide a single solution for
473 real-time, high-quality sequencing of long reads will be essential for implementing near real-time
474 diagnostics of infectious diseases and characterisation of pathogens.

475 **AUTHOR STATEMENTS**

476 **Authors and contributors**

477 The study was conceptualised by T.V.S. and U.S.J. Funding was secured by T.V.S., M.K., H.H. and U.S.J. Data
478 curation and investigation was performed by T.V.S.. Formal analysis was done by T.V.S. and S.O-P..
479 Resources were provided by M.K., T.V.S., H.H. and U.S.J.. U.S.J, H.H. and M.K. supervised the work. T.V.S
480 wrote the original draft and edited the manuscript. U.S.J., M.K., T.V.S., H.H., S.O-P. and H.M.W revised the
481 manuscript.

482 **Conflicts of interest**

483 The authors declare that there are no conflicts of interest

484 **Funding information**

490 **Ethical approval**

491 Isolates were obtained as part of routine clinical care, and details about the isolates have previously been
492 published. No ethical approvals were required.

493 **Acknowledgements**

## References

500

501  1.  Wexler HM. Bacteroides: the good, the bad, and the nitty-gritty. Clin Microbiol Rev. 2007
502      Oct;20(4):593–621. doi: 10.1128/CMR.00008-07

503  2.  Nagy E, Urbán E, Nord CE. Antimicrobial susceptibility of Bacteroides fragilis group isolates in Europe:
504      20 years of experience. Clin Microbiol Infect Off Publ Eur Soc Clin Microbiol Infect Dis. 2011
505      Mar;17(3):371–9. doi: 10.1111/j.1469-0691.2010.03256.x

506  3.  Ferløv-Schwensen SA, Sydenham TV, Hansen KC, Hoegh SV, Justesen US. Prevalences of antimicrobial
507      resistance and the cfiA resistance gene in danish bacteroides fragilis group isolates since 1973. Int J
508      Antimicrob Agents. 2017 Jun 27; doi: 10.1016/j.ijantimicag.2017.05.007

509  4.  Nagy E, Justesen US, Eitel Z, Urbán E. Development of EUCAST disk diffusion method for susceptibility
510      testing of the Bacteroides fragilis group isolates. Anaerobe. 2015 Feb;31:65–71. doi:
511      10.1016/j.anaerobe.2014.10.008

512  5.  Zankari E, Hasman H, Kaas RS, Seyfarth AM, Agersø Y, Lund O, et al. Genotyping using whole-genome
513      sequencing is a realistic alternative to surveillance based on phenotypic antimicrobial susceptibility
514      testing. J Antimicrob Chemother. 2013 Apr;68(4):771–7. doi: 10.1093/jac/dks496

515  6.  Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH, Del Ojo Elias C, et al. Predicting antimicrobial
516      susceptibilities for Escherichia coli and Klebsiella pneumoniae isolates using whole genomic sequence
517      data. J Antimicrob Chemother [Internet]. 2013 May 30 [cited 2013 Aug 19]; doi: 10.1093/jac/dkt180

518  7.  Boolchandani M, D'Souza AW, Dantas G. Sequencing-based methods and resources to study
519      antimicrobial resistance. Nat Rev Genet. 2019 Mar 18;1. doi: 10.1038/s41576-019-0108-4

520  8.  Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. Transforming clinical microbiology with
521      bacterial genome sequencing. Nat Rev Genet. 2012 Sep;13(9):601–12. doi: 10.1038/nrg3226

522  9.  Davies TJ, Stoesser N, Sheppard AE, Abuoun M, Fowler PW, Swann J, et al. Reconciling the potentially
523      irreconcilable? Genotypic and phenotypic amoxicillin-clavulanate resistance in Escherichia coli.
524      bioRxiv. 2019 Jan 7;511402. doi: 10.1101/511402

525  10. Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, et al. The role of whole
526      genome sequencing in antimicrobial susceptibility testing of bacteria: report from the EUCAST
527      Subcommittee. Clin Microbiol Infect Off Publ Eur Soc Clin Microbiol Infect Dis. 2017 Jan;23(1):2–22.
528      doi: 10.1016/j.cmi.2016.11.012

529  11. Nagy E, Becker S, Sóki J, Urbán E, Kostrzewa M. Differentiation of division I (cfiA-negative) and division
530      II (cfiA-positive) Bacteroides fragilis strains by matrix-assisted laser desorption/ionization time-of-
531      flight mass spectrometry. J Med Microbiol. 2011 Nov;60(Pt 11):1584–90. doi: 10.1099/jmm.0.031336-
532      0

533  12. Rogers MB, Parker AC, Smith CJ. Cloning and characterization of the endogenous cephalosporinase
534      gene, cepA, from Bacteroides fragilis reveals a new subgroup of Ambler class A beta-lactamases.
535      Antimicrob Agents Chemother. 1993 Nov;37(11):2391–400. doi: 10.1128/AAC.37.11.2391

536  13. Rasmussen BA, Gluzman Y, Tally FP. Cloning and sequencing of the class B beta-lactamase gene (ccrA)
537      from Bacteroides fragilis TAL3636. Antimicrob Agents Chemother. 1990 Aug;34(8):1590–2.

14

538  14.  Sydenham TV, Sóki J, Hasman H, Wang M, Justesen US, ESGAI (ESCMID Study Group on Anaerobic
539       Infections). Identification of antimicrobial resistance genes in multidrug-resistant clinical Bacteroides
540       fragilis isolates by whole genome shotgun sequencing. Anaerobe. 2015 Feb;31:59–64. doi:
541       10.1016/j.anaerobe.2014.10.009

542  15.  Ricker N, Qian H, Fulthorpe RR. The limitations of draft assemblies for understanding prokaryotic
543       adaptation and evolution. Genomics. 2012 Sep;100(3):167–75. doi: 10.1016/j.ygeno.2012.06.009

544  16.  Page AJ, De Silva N, Hunt M, Quail MA, Parkhill J, Harris SR, et al. Robust high-throughput prokaryote
545       de novo assembly and improvement pipeline for Illumina data. Microb Genomics [Internet]. 2016 Aug
546       25 [cited 2019 Apr 17];2(8). doi: 10.1099/mgen.0.000083

547  17.  Schmid M, Frei D, Patrignani A, Schlapbach R, Frey JE, Remus-Emsermann MNP, et al. Pushing the
548       limits of de novo genome assembly for complex prokaryotic genomes harboring very long, near
549       identical repeats. Nucleic Acids Res. 2018 Sep 28;46(17):8953–65. doi: 10.1093/nar/gky726

550  18.  Ring N, Abrahams JS, Jain M, Olsen H, Preston A, Bagby S. Resolving the complex Bordetella pertussis
551       genome using barcoded nanopore sequencing. Microb Genomics. 2018 Nov 21; doi:
552       10.1099/mgen.0.000234

553  19.  Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with multiplex MinION
554       sequencing. Microb Genomics [Internet]. 2017 [cited 2018 Apr 3];3(10). doi: 10.1099/mgen.0.000132

555  20.  Maio ND, Shaw LP, Hubbard A, George S, Sanderson N, Swann J, et al. Comparison of long-read
556       sequencing technologies in the hybrid assembly of complex bacterial genomes. bioRxiv. 2019 Jan
557       28;530824. doi: 10.1101/530824

558  21.  Cerdeño-Tárraga AM, Patrick S, Crossman LC, Blakely G, Abratt V, Lennard N, et al. Extensive DNA
559       inversions in the B. fragilis genome control variable gene expression. Science. 2005 Mar
560       4;307(5714):1463–5. doi: 10.1126/science.1107008

561  22.  Kuwahara T, Yamashita A, Hirakawa H, Nakayama H, Toh H, Okada N, et al. Genomic analysis of
562       Bacteroides fragilis reveals extensive DNA inversions regulating cell surface adaptation. Proc Natl
563       Acad Sci U S A. 2004 Oct 12;101(41):14919–24. doi: 10.1073/pnas.0404172101

564  23.  Patrick S, Blakely GW, Houston S, Moore J, Abratt VR, Bertalan M, et al. Twenty-eight divergent
565       polysaccharide loci specifying within- and amongst-strain capsule diversity in three strains of
566       Bacteroides fragilis. Microbiology. 2010 Nov;156(Pt 11):3255–69. doi: 10.1099/mic.0.042978-0

567  24.  Nikitina AS, Kharlampieva DD, Babenko VV, Shirokov DA, Vakhitova MT, Manolov AI, et al. Complete
568       Genome Sequence of an Enterotoxigenic Bacteroides fragilis Clinical Isolate. Genome Announc
569       [Internet]. 2015 May 7;3(3). doi: 10.1128/genomeA.00450-15

570  25.  Risse J, Thomson M, Patrick S, Blakely G, Koutsovoulos G, Blaxter M, et al. A single chromosome
571       assembly of Bacteroides fragilis strain BE1 from Illumina and MinION nanopore sequencing data.
572       GigaScience [Internet]. 2015 Dec 4 [cited 2017 Jun 22];4. doi: 10.1186/s13742-015-0101-6

573  26.  Soki, Jozsef. Bacteroides fragilis S14 genome sequencing and assembly (Data accessed on NCBI RefSeq
574       database accession GCF_001682215.1). 2015.

575   27.   Ho P-L, Yau C-Y, Wang Y, Chow K-H. Determination of the mutant–prevention concentration of
576         imipenem for the two imipenem–susceptible Bacteroides fragilis strains, Q1F2 (cfiA-positive) and
577         ATCC 25282 (cfiA-negative). Int J Antimicrob Agents. 2018 Feb 1;51(2):270–1. doi:
578         10.1016/j.ijantimicag.2017.08.004

579   28.   Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. Bioinformatics. 2018
580         Sep 1;34(17):i884–90. doi: 10.1093/bioinformatics/bty560

581   29.   Desai A, Marwah VS, Yadav A, Jha V, Dhaygude K, Bangar U, et al. Identification of Optimum
582         Sequencing Depth Especially for De Novo Genome Assembly of Small Genomes Using Next Generation
583         Sequencing Data. PLoS ONE [Internet]. 2013 Apr 12 [cited 2014 May 22];8(4). doi:
584         10.1371/journal.pone.0060204

585   30.   Wick RR, Judd LM, Holt KE. Deepbinner: Demultiplexing barcoded Oxford Nanopore reads with deep
586         convolutional neural networks. bioRxiv. 2018 Sep 14;366526. doi: 10.1101/366526

587   31.   De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack: visualizing and processing
588         long-read sequencing data. Bioinformatics. 2018 Aug 1;34(15):2666–9. doi:
589         10.1093/bioinformatics/bty149

590   32.   Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies.
591         Bioinformatics. 2013 Apr 15;29(8):1072–5. doi: 10.1093/bioinformatics/btt086

592   33.   Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of
593         microbial genomes recovered from isolates, single cells, and metagenomes. Genome Res. 2015 Jul
594         1;25(7):1043–55. doi: 10.1101/gr.186072.114

595   34.   Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al. BUSCO
596         Applications from Quality Assessments to Gene Prediction and Phylogenomics. Mol Biol Evol. 2018
597         Mar 1;35(3):543–8. doi: 10.1093/molbev/msx319

598   35.   Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics. 2014 Jul 15;30(14):2068–9.
599         doi: 10.1093/bioinformatics/btu153

600   36.   Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition.
601         Proc Natl Acad Sci. 2009 Oct 22;pnas.0906412106. doi: 10.1073/pnas.0906412106

602   37.   Clark SC, Egan R, Frazier PI, Wang Z. ALE: a generic assembly likelihood evaluation framework for
603         assessing the accuracy of genome and metagenome assemblies. Bioinformatics. 2013 Feb
604         15;29(4):435–43. doi: 10.1093/bioinformatics/bts723

605   38.   Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo genome
606         assemblies. Bioinforma Oxf Engl. 2015 Oct 15;31(20):3350–2. doi: 10.1093/bioinformatics/btv383

607   39.   Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 2018 Sep
608         15;34(18):3094–100. doi: 10.1093/bioinformatics/bty191

609   40.   Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinforma
610         Oxf Engl. 2009 Jul 15;25(14):1754–60. doi: 10.1093/bioinformatics/btp324

611   41.   Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA, Harris SR. Circlator: automated circularization of
612         genome assemblies using long sequencing reads. Genome Biol. 2015 Dec 29;16:294. doi:
613         10.1186/s13059-015-0849-0

614   42.   Tatusova T, DiCuccio M, Badretdin A, Chetvernin V, Nawrocki EP, Zaslavsky L, et al. NCBI prokaryotic
615         genome annotation pipeline. Nucleic Acids Res. 2016 19;44(14):6614–24. doi: 10.1093/nar/gkw569

616   43.   Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, et al. Identification of
617         acquired antimicrobial resistance genes. J Antimicrob Chemother. 2012 Nov;67(11):2640–4. doi:
618         10.1093/jac/dks261

619   44.   Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P, Tsang KK, et al. CARD 2017: expansion and model-
620         centric curation of the comprehensive antibiotic resistance database. Nucleic Acids Res. 2017
621         04;45(D1):D566–73. doi: 10.1093/nar/gkw1004

622   45.   Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. ISfinder: the reference centre for bacterial
623         insertion sequences. Nucleic Acids Res. 2006 Jan 1;34(Database issue):D32-36. doi:
624         10.1093/nar/gkj014

625   46.   Galata V, Fehlmann T, Backes C, Keller A. PLSDB: a resource of complete bacterial plasmids. Nucleic
626         Acids Res. 2019 Jan 8;47(D1):D195–202. doi: 10.1093/nar/gky1050

627   47.   Jørgensen TS, Xu Z, Hansen MA, Sørensen SJ, Hansen LH. Hundreds of Circular Novel Plasmids and
628         DNA Elements Identified in a Rat Cecum Metamobilome. PLoS ONE [Internet]. 2014 Feb 4 [cited 2019
629         Mar 12];9(2). doi: 10.1371/journal.pone.0087924

630   48.   Shah HN. The Genus Bacteroides and Related Taxa. In: Balows A, Trüper HG, Dworkin M, Harder W,
631         Schleifer K-H, editors. The Prokaryotes: A Handbook on the Biology of Bacteria: Ecophysiology,
632         Isolation, Identification, Applications [Internet]. New York, NY: Springer New York; 1992 [cited 2019
633         Mar 21]. p. 3593–607. doi: 10.1007/978-1-4757-2191-1_34

634   49.   Sóki J, Wareham DW, Rátkai C, Aduse-Opoku J, Urbán E, Nagy E. Prevalence, nucleotide sequence and
635         expression studies of two proteins of a 5.6kb, class III, Bacteroides plasmid frequently found in clinical
636         isolates from European countries. Plasmid. 2010 Mar;63(2):86–97. doi:
637         10.1016/j.plasmid.2009.12.002

638   50.   Database resources of the National Center for Biotechnology Information. Nucleic Acids Res. 2016 Jan
639         4;44(Database issue):D7–19. doi: 10.1093/nar/gkv1290

640   51.   Trinh S, Haggoud A, Reysset G, Sebald M. Plasmids pIP419 and pIP421 from Bacteroides: 5-
641         nitroimidazole resistance genes and their upstream insertion sequence elements. Microbiol Read
642         Engl. 1995 Apr;141 ( Pt 4):927–35. doi: 10.1099/13500872-141-4-927

643   52.   Haggoud A, Trinh S, Moumni M, Reysset G. Genetic analysis of the minimal replicon of plasmid pIP417
644         and comparison with the other encoding 5-nitroimidazole resistance plasmids from Bacteroides spp.
645         Plasmid. 1995 Sep;34(2):132–43. doi: 10.1006/plas.1995.9994

646   53.   Sóki J, Gal M, Brazier JS, Rotimi VO, Urbán E, Nagy E, et al. Molecular investigation of genetic
647         elements contributing to metronidazole resistance in Bacteroides strains. J Antimicrob Chemother.
648         2006 Feb;57(2):212–20. doi: 10.1093/jac/dki443

649   54.   Hartmeyer GN, Sóki J, Nagy E, Justesen US. Multidrug-resistant Bacteroides fragilis group on the rise
650         in Europe? J Med Microbiol. 2012 Dec;61(Pt 12):1784–8. doi: 10.1099/jmm.0.049825-0

651   55.   Darling AE, Mau B, Perna NT. progressiveMauve: Multiple Genome Alignment with Gene Gain, Loss
652         and Rearrangement. PLOS ONE. 2010 Jun 25;5(6):e11147. doi: 10.1371/journal.pone.0011147

653   56.   Nishida H. Comparative Analyses of Base Compositions, DNA Sizes, and Dinucleotide Frequency
654         Profiles in Archaeal and Bacterial Chromosomes and Plasmids [Internet]. International Journal of
655         Evolutionary Biology. 2012 [cited 2019 Apr 26]. doi: 10.1155/2012/342482

656   57.   Brettin T, Davis JJ, Disz T, Edwards RA, Gerdes S, Olsen GJ, et al. RASTtk: a modular and extensible
657         implementation of the RAST algorithm for building custom annotation pipelines and annotating
658         batches of genomes. Sci Rep. 2015 Feb 10;5:8365. doi: 10.1038/srep08365

659   58.   Bayley DP, Rocha ER, Smith CJ. Analysis of cepA and other Bacteroides fragilis genes reveals a unique
660         promoter structure. FEMS Microbiol Lett. 2000 Dec 1;193(1):149–54. doi: 10.1111/j.1574-
661         6968.2000.tb09417.x

662   59.   Lin Y, Yuan J, Kolmogorov M, Shen MW, Chaisson M, Pevzner PA. Assembly of long error-prone reads
663         using de Bruijn graphs. Proc Natl Acad Sci. 2016 Dec 27;113(52):E8396–405. doi:
664         10.1073/pnas.1604560113

665   60.   Kamath GM, Shomorony I, Xia F, Courtade TA, Tse DN. HINGE: long-read assembly achieves optimal
666         repeat resolution. Genome Res. 2017;27(5):747–56. doi: 10.1101/gr.216465.116

667   61.   Payne A, Holmes N, Rakyan V, Loose M. Whale watching with BulkVis: A graphical viewer for Oxford
668         Nanopore bulk fast5 files. bioRxiv. 2018 May 3;312256. doi: 10.1101/312256

669   62.   de Lannoy C, Risse J, de Ridder D. poreTally: run and publish de novo nanopore assembler
670         benchmarks. Bioinformatics [Internet]. [cited 2019 Apr 26]; doi: 10.1093/bioinformatics/bty1045

671   63.   Carattoli A, Zankari E, García-Fernández A, Larsen MV, Lund O, Villa L, et al. In Silico Detection and
672         Typing of Plasmids using PlasmidFinder and Plasmid Multilocus Sequence Typing. Antimicrob Agents
673         Chemother. 2014 Jul 1;58(7):3895–903. doi: 10.1128/AAC.02412-14

674   64.   Arredondo-Alonso S, Rogers MRC, Braat JC, Verschuuren TD, Top J, Corander J, et al. mlplasmids: a
675         user-friendly tool to predict plasmid- and chromosome-derived sequences for single species. Microb
676         Genomics [Internet]. 2018 [cited 2019 May 9];4(11). doi: 10.1099/mgen.0.000224

677   65.   Nguyen M, Vedantam G. Mobile genetic elements in the genus Bacteroides, and their mechanism(s)
678         of dissemination. Mob Genet Elem. 2011;1(3):187–96. doi: 10.4161/mge.1.3.18448

679   66.   Shkoporov AN, Khokhlova EV, Kulagina EV, Smeianov VV, Kuchmiy AA, Kafarskaya LI, et al. Analysis of
680         a novel 8.9kb cryptic plasmid from Bacteroides uniformis, its long-term stability and spread within
681         human microbiota. Plasmid. 2013 Mar 1;69(2):146–59. doi: 10.1016/j.plasmid.2012.11.002

682   67.   McNulty NP, Wu M, Erickson AR, Pan C, Erickson BK, Martens EC, et al. Effects of diet on resource
683         utilization by a model human gut microbiota containing Bacteroides cellulosilyticus WH2, a symbiont
684         with an extensive glycobiome. PLoS Biol. 2013;11(8):e1001637. doi: 10.1371/journal.pbio.1001637

685    68.    Pierce JV, Bernstein HD. Genomic Diversity of Enterotoxigenic Strains of Bacteroides fragilis. PLoS ONE
686           [Internet]. 2016 Jun 27 [cited 2019 Mar 26];11(6). doi: 10.1371/journal.pone.0158171

687    69.    Husain F, Veeranagouda Y, Hsi J, Meggersee R, Abratt V, Wexler HM. Two multidrug-resistant clinical
688           isolates of Bacteroides fragilis carry a novel metronidazole resistance nim gene (nimJ). Antimicrob
689           Agents Chemother. 2013 Aug;57(8):3767–74. doi: 10.1128/AAC.00386-13

690    70.    Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. bioRxiv. 2019 Jan 26;530972. doi:
691           10.1101/530972

692    71.    Li H. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences.
693           Bioinformatics. 2016 Jul 15;32(14):2103–10. doi: 10.1093/bioinformatics/btw152

694    72.    Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs.
695           Nat Biotechnol. 2019 Apr 1;1. doi: 10.1038/s41587-019-0072-8

696    73.    Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-
697           read assembly via adaptive k-mer weighting and repeat separation. Genome Res. 2017 May
698           1;27(5):722–36. doi: 10.1101/gr.215087.116

699    74.    Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A New Genome
700           Assembly Algorithm and Its Applications to Single-Cell Sequencing. J Comput Biol. 2012
701           May;19(5):455–77. doi: 10.1089/cmb.2012.0021

702    75.    Antipov D, Korobeynikov A, McLean JS, Pevzner PA. hybridSPAdes: an algorithm for hybrid assembly
703           of short and long reads. Bioinformatics. 2016 Apr 1;32(7):1009–15. doi:
704           10.1093/bioinformatics/btv688

705    76.    Souvorov A, Agarwala R, Lipman DJ. SKESA: strategic k-mer extension for scrupulous assemblies.
706           Genome Biol. 2018 Oct 4;19(1):153. doi: 10.1186/s13059-018-1540-z

707    77.    Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies from short
708           and long sequencing reads. PLOS Comput Biol. 2017 Aug 6;13(6):e1005595. doi:
709           10.1371/journal.pcbi.1005595

710    78.    Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled *de novo* using only nanopore
711           sequencing data. Nat Methods. 2015 Aug;12(8):733–5. doi: 10.1038/nmeth.3444

712    79.    Vaser R, Sovic I, Nagarajan N, Sikic M. Fast and accurate de novo genome assembly from long
713           uncorrected reads. Genome Res. 2017 Jan 18;gr.214270.116. doi: 10.1101/gr.214270.116

714    80.    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An Integrated Tool for
715           Comprehensive Microbial Variant Detection and Genome Assembly Improvement. PLOS ONE. 2014
716           Nov 19;9(11):e112963. doi: 10.1371/journal.pone.0112963

717    81.    Krumsiek J, Arnold R, Rattei T. Gepard: a rapid and sensitive tool for creating dotplots on genome
718           scale. Bioinforma Oxf Engl. 2007 Apr 15;23(8):1026–8. doi: 10.1093/bioinformatics/btm039

719    82.    Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer. Bioinformatics. 2011 Apr
720           1;27(7):1009–10. doi: 10.1093/bioinformatics/btr039

721    83.  Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better, faster version of the
722         PHAST phage search tool. Nucleic Acids Res. 2016 08;44(W1):W16-21. doi: 10.1093/nar/gkw387

723

724 **TABLES AND FIGURES**

725

| Genome assembler and version | Link | Reference |
|---|---|---|
| Wtdbg2 v2.3 | https://github.com/ruanjue/wtdbg2 | (70) |
| Miniasm v0.3r179 | https://github.com/lh3/miniasm | (39,71) |
| Flye v2.3.7 | https://github.com/fenderglass/Flye | (59,72) |
| Canu v1.8 | https://github.com/marbl/canu | (73) |
| Spades (including Hybridspades) v3.13.0 | https://github.com/ablab/spades | (74,75) |
| Skesa v2.3.0 | https://github.com/ncbi/SKESA | (76) |
| Unicycler v0.4.7 | https://github.com/rrwick/Unicycler | (77) |
| **Assembly polishing tools** | | |
| Nanopolish v0.10.2 | https://github.com/jts/nanopolish | (78) |
| Racon v1.3.1 | https://github.com/isovic/racon | (79) |
| Pilon v1.22 | https://github.com/broadinstitute/pilon | (80) |

726     Table 1 - Genome assemblers and polishing tools tested

21

727

| Assembly | Contigs | Largest contig | Total length | Mis-assemblies | Genome fraction (%) | Mismatches per 100 kbp | Indels per 100 kbp | Average nucleotide identity | CheckM Completeness | BUSCOs: complete and single-copy/ complete and duplicate/ fragment (of 443) | Prokka genes | Prokka rRNA | Prokka tRNA | Total ALE score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GCF_000025985.1 | 2 | 5,205,140 | 5,241,700 | 0 | 100.000 | 0 | 0 | 100.0000 | 99.26 | 442/0/1 | 4439 | 19 | 73 | -17071758.95 |
| Skesa | 46 | 553,341 | 5,201,945 | 3 | 99.237 | 0.23 | 0.15 | 99.9980 | 99.26 | 440/2/1 | 4391 | 2 | 62 | -20926329.69 |
| Spades | 23 | 1,779,941 | 5,212,217 | 4 | 99.396 | 0.44 | 0.17 | 99.9866 | 99.26 | 440/2/1 | 4407 | 3 | 56 | -19676529.39 |
| Canu.OF.CO.RO2.RI.PI3 | 2 | 5,247,938 | 5,350,432 | 8 | 99.972 | 4.94 | 15.9 | 99.9752 | 99.26 | 442/0/1 | 4634 | 19 | 73 | -19283611.73 |
| Flye.OF.CS.PI5.RI | 5 | 2,282,650 | 5,269,269 | 4 | 99.917 | 1.07 | 6.24 | 99.9781 | 99.26 | 441/1/1 | 4476 | 19 | 73 | -18222322.23 |
| Miniasm.OF.CM.RO2.PI5 | 3 | 5,204,445 | 5,277,434 | 2 | 99.972 | 5.21 | 17.75 | 99.9691 | 98.88 | 442/0/1 | 4607 | 19 | 73 | -17789234.97 |
| Wtdbg2.OF.CO.RO2.PI6.RI | 3 | 5,192,352 | 5,234,448 | 7 | 99.723 | 3.23 | 3.04 | 99.9807 | 99.26 | 442/0/1 | 4437 | 19 | 73 | -18750266.21 |
| SpadesHybrid.CS | 5 | 3,093,122 | 5,242,724 | 7 | 99.987 | 1.89 | 0.53 | 99.9856 | 99.26 | 440/2/1 | 4441 | 19 | 73 | -18535980.68 |
| Unicycler.OF.CS | 2 | 5,205,133 | 5,241,693 | 2 | 99.972 | 0.84 | 0.48 | 99.9997 | 99.26 | 442/0/1 | 4435 | 19 | 73 | -17200232.52 |

728 Table 2. Selected quality indicators for the best genome assembly of *B. fragilis* CCUG4856T per assembly pipeline. RefSeq accession GCF_000025985.1
729 was used as reference. OF: ONT reads filtered with Filtlong, CS: Canu corrected standard settings, CM: Canu corrected with option corMinCoverage=0,
730 CO: Canu corrected with option coroutCoverage=999, RO2: Two rounds of Racon polishing with ONT reads, RI: Racon polishing with Illumina reads,
731 PI[n]: Pilon polishing with Illumina reads, [n] rounds. Full results are available in Supplementary Table S2.

732

733

| Assembly | Total length (bp) | Largest contig (bp) | Local mis-assembl ies | Genome fraction (%) | Mismatch es per 100 kbp | Indels per 100 kbp | K-mer-based compl. (%) | K-mer-based misjoin s | Average nucletide identity | Prokka CDS | Prokka genes | Total ALE score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GCF_00002 5985.1 | 5,241,700 | 5,205,140 | 0 | 100.000 | 0 | 0 | 100.00 | 0 | 100.00000 | 4346 | 4439 | -17071758.95 |
| OF | 5,241,602 | 5,205,042 | 3 | 99.970 | 1.11 | 0.65 | 99.96 | 0 | 99.99896 | 4343 | 4436 | -17245134.52 |
| OF.RI | 5,241,606 | 5,205,046 | 3 | 99.970 | 1.09 | 0.67 | 99.96 | 3 | 99.99887 | 4345 | 4438 | -17247815.86 |
| OF.CS | 5,241,693 | 5,205,133 | 2 | 99.972 | 0.84 | 0.48 | 99.97 | 1 | 99.99974 | 4342 | 4435 | -17200232.52 |
| OF.CS.RI | 5,241,698 | 5,205,138 | 2 | 99.972 | 0.88 | 0.52 | 99.96 | 1 | 99.99968 | 4346 | 4439 | -17206271.66 |
| OF.CM | 5,241,691 | 5,205,131 | 2 | 99.972 | 0.88 | 0.5 | 99.96 | 1 | 99.99966 | 4343 | 4436 | -17201292.44 |
| OF.CM.RI | 5,241,696 | 5,205,136 | 2 | 99.972 | 0.95 | 0.55 | 99.97 | 1 | 99.99975 | 4343 | 4436 | -17193184.79 |
| OF.CO | 5,241,693 | 5,205,133 | 2 | 99.972 | 0.84 | 0.48 | 99.97 | 1 | 99.99974 | 4342 | 4435 | -17200232.52 |
| OF.CO.RI | 5,241,698 | 5,205,138 | 2 | 99.972 | 0.88 | 0.52 | 99.96 | 1 | 99.99968 | 4346 | 4439 | -17206271.66 |

734 Table 3. Hybrid Unicycler assemblies of *B. fragilis* CCUG4856T. RefSeq accession GCF_000025985.1 was used as reference. OF: ONT reads filtered with
735 Filtlong; CS: Canu corrected standard settings; CM: Canu corrected with option corMinCoverage=0; CO: Canu corrected with option
736 coroutCoverage=999; RI: Racon polishing with Illumina reads. Unicycler performs assembly polishing with Racon (ONT reads) and Pilon. Full results are
737 available in Supplementary Table S2.

738

739

740 Figure 1. Dot plot matrix of the alignment of the reference assembly and the hybrid Unicycler assembly
741 using Gepard v1.40 (81). The *B. fragilis* NCTC9343 (RefSeq GCF_000025985.1) reference assembly derived
742 from Sanger sequencing is on the x-axis and the hybrid Unicycler assembly on the y-axis. On this otherwise
743 near perfect alignment with high similarity, an 88,045 bp inversion with 100% ID is observed at nucleotide
744 positions 2,941,962..3,030,006 on the Unicycler assembly (2,005,742..2,093,786 on the reference
745 sequence) (indicated by the blue arrow).

24

| Strain | Best Spades graph | Hybrid assembly after first nanopore sequencing run. No Canu correction | Hybrid assembly after first nanopore sequencing run with Canu correction | Hybrid assembly with data from supplementary runs with Canu correction | After manual finishing | Final assemblies contig statistics | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | No. | Length | Relative read depth |
| CCUG 4856T | | | | | | 1 | 5,205,133 | 1.00x |
| | | | | | | 2 | 36,560 | 7.42x |
| BFO17 | | | | | | 1 | 5,474,541 | 1.00x |
| | | | | | | 2 | 85,671 | 1.85x |
| | | | | | | 3 | 5,594 | 23.03x |
| BFO18 | | | | | | 1 | 5,302,644 | 1.00x |
| | | | | | | 2 | 7,221 | 25.98x |
| | | | | | | 3 | 4,137 | 50.80x |
| | | | | | | 4 | 2,782 | 59.91x |
| S01 | | | | | | 1 | 5,325,251 | 1.00x |
| | | | | | | 2 | 78,085 | 2.29x |
| | | | | | | 3 | 8,331 | 20.99x |
| | | | | | | 4 | 5,595 | 22.67x |
| BFO42 | | | | | | 1 | 5,141,257 | 1.00x |
| | | | | | | 2 | 8,306 | 40.06x |
| | | | | | | 3 | 5,594 | 40.15x |
| BFO67 | | | | | | 1 | 5,478,614 | 1.00x |
| | | | | | | 2 | 6,129 | 94.15x |
| BFO85 | | | | | | 1 | 5,504,076 | 1.00x |

746

Figure 2. Evolution of genome assemblies with added data and manual finishing. The best SPAdes assembly graphs by Unicycler with short reads only are shown on the far left. Supplying ONT reads improved the assemblies overall, but only three were circularised with singular chromosome contigs with data from the initial MinION sequencing runs. Adding additional ONT data and correcting reads with Canu did not improve assemblies for all isolates. Manual finishing was necessary to finish assemblies for three isolates. Assembly graph images generated with Bandage. Read information can be found in Supplementary Table S1.

| Strain | Sequence | Length (bp) | Relative read depth | GC% | PLSDB results | | | | Plasmid replicon family (%COV, %ID) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | Best hit accession no. | Plasmid hit name | %ID | Length of the sequence of best hit (bp) | |
| CCUG4856T | Chr. | 5,205,133 | 1.00x | 43.19 | - | - | - | - | - |
| | pBF9343 | 36,560 | 7.42x | 32.19 | NC_006873.1 | pBF9343 | 100 | 36,560 | Rep_3 (100/100) |
| BFO17 | Chr. | 5,474,541 | 1.00x | 43.51 | - | - | - | - | - |
| | pBFO17_1 | 85,671 | 1.85x | 36.78 | NC_006873.1 | pBF9343 | 80.7 | 36,560 | none |
| | pBFO17_2 | 5,594 | 23.03x | 39.65 | NC_011073.1 | pBFP35 | 99.9 | 5,594 | Rep_1 (100/100) |
| BFO18 | Chr. | 5,302,644 | 1.00x | 43.34 | - | - | - | - | - |
| | pBFO18_1 | 7,221 | 25.98x | 42.32 | NC_015168.1 | pBACSA02 | 85.6 | 19,280 | Rep_3 (99.69/99.69) |
| | pBFO18_2 | 4,137 | 50.80x | 45.40 | NC_019534.1 | pBFUK1 | 92.2 | 12,817 | Rep_3 (100.00/98.24) |
| | pBFO18_3 | 2,782 | 59.91x | 41.45 | NC_005026.1 | pBI143 | 94.6 | 2,747 | RepL (89.66/49.22)[a] |
| S01 | Chr. | 5,325,251 | 1.00x | 43.57 | - | - | - | - | - |
| | pBFS01_1 | 78,085 | 2.29x | 36.04 | NC_006873.1 | pBF9343 | 80.7 | 36,560 | none |
| | pBFS01_2 | 8,331 | 20.99x | 41.17 | NC_015166.1 | pBACSA03 | 95.6 | 6,277 | Rep_3 (100.00/97.85) |
| | pBFS01_3 | 5,595 | 22.67x | 39.62 | NC_011073.1 | pBFP35 | 99.9 | 5,594 | Rep_1 (100.00/99.48) |
| BFO42 | Chr. | 5,141,257 | 1.00x | 43.35 | - | - | - | - | - |
| | pBFO32_1 | 8,306 | 40.06x | 43.34 | KJ830768.1 | pBF69566b | 96.0 | 11,019 | RHH_1 (92.94/64.63) Rep_3 (93.64/68.31) |
| | pBFO32_2 | 5,594 | 40.15x | 39.63 | NC_011073.1 | pBFP35 | 99.9 | 5,594 | Rep_1 (100.00 /99.48) |
| BFO67 | Chr. | 5,478,614 | 1.00x | 43.85 | - | - | - | - | - |
| | pBFO67_1 | 6,129 | 94.15x | 41.67 | NC_011073.1 | pBFP35 | 76.9 | 5,594 | Rep_3 (100.00/99.69) |
| BFO85 | Chr. | 5,504,076 | 1.00x | 43.60 | - | - | - | - | - |

754 Table 4. Putative plasmid sequences of the complete *B. fragilis* assemblies. Putative plasmid sequences from the hybrid assemblies of *B. fragilis*
755 CCUG4856T and the six MDR *B. fragilis* isolates were screened using the PLSDB. The best hit to plasmids from cultured isolates is shown. Only three
756 putative plasmids from the MDR *B. fragilis* isolate assemblies could be identified with confident %ID. For most sequences, plasmid replication family
757 proteins were identified in the putative plasmids using ABRicate with a database of sequences downloaded from the Pfam database, strengthening

758 the interpretation that the circularised putative plasmid sequences do in fact represent plasmids harboured by the isolates. Notes: [a]Annotated as
759 RepA protein in the PGAP annotation. Abbreviations: %ID, %COV, no.; number, Chr.; chromosome.
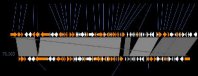
760

Figure 3. Linear representation of an alignment of putative circular plasmid sequences pBFO17_1 pBFS01_1. Comparison of the putative circular plasmids pBFO17_1 and pBFS01_1 (reverse complement for better visualisation) using EasyFig (82). EasyFig uses BLAST to identify sequences of similarity. Sequence similarities of >98% is indicated by full colouring, a darker colour indicates a higher %ID. Products of annotated CDS' are shown. CDS' annotated as hypothetical or Domain of Unknown Function are coloured white. The two sequences show a very high degree of similarity. pBFO17_1 is 7,586 bp longer than pBFS01_1. This is mainly due to the insertion of a reverse transcriptase (pBFO17_1, 11367..13034) (disrupting a DNA methylase), the insertion of prophage from position 56125 to 61162) (identified as an incomplete prophage using PHASTER (83)) and an IS*1380* family-like transposase (67933..69237). The regions pBFO17_1 50711..52501 and pBFS01_1 32248..30304 are not similar. Possibly, the insertion of two transposases in pBFO17_1 have excised most of the ParB-family DNA partitioning protein in the corresponding sequence range in pBFS01_1.

28

| Antimicrobial susceptibility[a] | | | | Antimicrobial resistance genes and IS elements[b] | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Strain | Antimi-crobial | Etest MIC (mg/L) | Result | Gene | Upstream IS element | Sequence[c] | %ID | %COV | Associated resistance to drug class |
| BFO17 | MEM | >32 | R | *cfiA11* | IS*614B* | Chr | 100.00 | 99.20 | Carbapenem |
| | IPM | >32 | R | | | | | | |
| | MTZ | >32 | R | *nimJ* | IS*614B* | Chr | 99.40 | 100.00 | Nitroimidazole |
| | | | | *nimJ* | IS*614B* | Chr$_c$ | 99.40 | 100.00 | Nitroimidazole |
| | CLI | 0.094 | S | | | | | | |
| | PTZ | >256 | R | | | | | | |
| | | | | *tetQ* | | Chr | 99.34 | 99.34 | Tetracycline |
| | | | | *cfxA4* | | Chr | 85.71 | 100.00 | Cephamycin |
| | | | | *bexB* | | Chr$_c$ | 91.21 | 100.00 | Fluoroquinolone |
| | | | | *bexA* | | Chr | 73.77 | 99.02 | Fluoroquinolone |
| BFO18 | MEM | >32 | R | *cfiA2_1* | IS*Bf12* | Chr | 100.00 | 100.00 | Carbapenem |
| | IPM | 16 | R | | | | 100.00 | 100.00 | |
| | MTZ | 16 | R | *nimD* | IS*1169* | 2 | 99.19 | 100.00 | Nitroimidazole |
| | CLI | 6 | R | *ermF[d]* | IS*4351* | Chr$_c$ | 99.83 | 72.03 | Clindamycin |
| | | | | | IS*Bthe1*[d] | Chr$_c$ | 70.97 | 97.19 | |
| | | | | *erm(F)[d]* | | Chr$_c$ | 99.58 | 29.71 | |
| | | | | *lnu(AN2)* | | Chr$_c$ | 100.00 | 100.00 | Clindamycin |
| | PTZ | >256 | R | | | | | | |
| | | | | *ugd* | | Chr | 65.69 | 53.04 | Polymyxin |
| | | | | *bexA* | | Chr$_c$ | 73.60 | 99.02 | Fluoroquinolone |
| | | | | *bexB* | | Chr | 91.14 | 100.00 | Fluoroquinolone |
| | | | | *tet(Q)* | | Chr$_c$ | 99.79 | 100.00 | Tetracycline |
| | | | | *mef(En2)* | | Chr$_c$ | 99.83 | 100.00 | Macrolides |
| S01 | MEM | >32 | R | *cfiA13_1* | IS*1187* | Chr | 99.20 | 100.00 | Carbapenem |
| | IPM | 16 | R | | | | | | |
| | MTZ | 64 | R | *nimE* | IS*Bf6* | 3 | 100.00 | 100.00 | Nitroimidazole |
| | CLI | >32 | R | *erm(F)* | IS*1187* | Chr | 99.50 | 100.00 | Clindamycin |
| | PTZ | 6 | S | | | | | | |
| | | | | *tetQ* | | Chr | 90.02 | 99.95 | Tetracycline |
| | | | | *tet(Q)* | | Chr$_c$ | 99.84 | 100.00 | Tetracycline |
| | | | | *bexB* | | Chr$_c$ | 91.06 | 100.00 | Fluoroquinolone |
| | | | | *bexA* | | Chr | 74.03 | 98.80 | Fluoroquinolone |
| BFO42 | MEM | 0.094 | S | | | | | | |
| | IPM | 0.25 | S | | | | | | |
| | MTZ | 8 | R | *nimA* | IS*Bf13* | 2 | 98.64 | 96.61 | Nitroimidazole |
| | CLI | >256 | R | *erm(F)* | IS*613* | Chr | 99.50 | 100.00 | Clindamycin |
| | | | | *lnu(AN2)* | | Chr | 100.00 | 100.00 | Clindamycin |
| | PTZ | 0.38 | S | | | | | | |

| Strain | Antimicrobial | MIC | S/R | Gene | IS element | Location | %ID | %COV | Resistance |
|---|---|---|---|---|---|---|---|---|---|
| | | | | ugd | | Chr | 70.38 | 31.45 | |
| | | | | cepA-49 | | Chr_c | 100.00 | 100.00 | Cephalosporin |
| | | | | mef(En2) | | Chr | 99.83 | 100.00 | Macrolide |
| | | | | ugd | | Chr | 71.15 | 31.11 | Polymyxin |
| | | | | tetQ | | Chr_c | 100.00 | 100.00 | Tetracycline |
| | | | | bexB | | Chr | 99.12 | 100.00 | Fluoroquinolone |
| | | | | ere(D) | | Chr | 96.66 | 100.00 | Erythromycin |
| | | | | aadS | | Chr_c | 99.88 | 100.00 | Aminoglycoside |
| | | | | OXA-347 | IS613 | Chr_c | 100.00 | 100.00 | Penicillin, cephalosporin |
| | | | | bexA | | Chr | 75.09 | 99.62 | Fluoroquinolone |
| BFO67 | MEM | 8 | R | cfiA13_1 | None | Chr | 100.00 | 100.00 | Carbapenem |
| | IPM | 0.5 | S | | | | | | |
| | MTZ | 0.19 | S | | | | | | |
| | CLI | 0.38 | S | | | | | | |
| | PTZ | 2 | S | | | | | | |
| | | | | cfxA2 | ISBf11 | Chr | 99.69 | 100.00 | Cephamycin |
| | | | | mef(En2) | | Chr | 99.75 | 100.00 | Macrolide |
| | | | | lnu(AN2) | | Chr | 100.00 | 100.00 | Clindamycin |
| | | | | ugd | | Chr_c | 66.76 | 56.30 | Polymyxin |
| | | | | tet(Q) | | Chr | 100.00 | 100.00 | Tetracycline |
| | | | | bexB | | Chr_c | 90.92 | 100.00 | Fluoroquinolone |
| | | | | bexA | | Chr | 73.90 | 99.02 | Fluoroquinolone |
| BFO85 | MEM | 32 | R | cfiA2_1 | None | Chr | 100.00 | 100.00 | Carbapenem |
| | IPM | 1 | S | | | | | | |
| | MTZ | 0.25 | S | | | | | | |
| | CLI | >256 | R | ermB | | Chr_c | 99.19 | 98.66 | Clindamycin |
| | PTZ | 2 | S | | | | | | |
| | | | | ugd | | Chr | 69.84 | 31.45 | Polymyxin |
| | | | | tetQ | | Chr_c | 90.02 | 99.95 | Tetracycline |
| | | | | aadE | | Chr_c | 100.00 | 100.00 | Aminoglycoside |
| | | | | aad9 | | Chr_c | 100.00 | 100.00 | Aminoglycoside |
| | | | | bexB | | Chr_c | 90.92 | 100.00 | Fluoroquinolone |
| | | | | bexA | | Chr | 73.53 | 99.02 | Fluoroquinolone |
| | | | | cfxA2 | IS614 | Chr_c | 100.00 | 100.00 | Cephamycin |
| | | | | tet(Q) | | Chr_c | 99.84 | 100.00 | Tetracycline |

Table 5. Antimicrobial susceptibility and resistance genes and IS elements for the six MDR *B. fragilis* strains. Identified genes are displayed next to the relevant antimicrobials. Identified IS elements in correct orientation (opposite strand) directly upstream of the genes are included. The %ID and %COV refer to the gene hit. Hits with %ID or %COV <98% were confirmed with BLASTx searches. The hits for *ugd* represent possible homologs for genes coding for PmrE, which is involved in polymyxin resistance in Gram-negative bacteria. Full ABRicate results with nucleotide positions and information on the IS elements is available the Supplementary Tables S4 AMR-IS-results.xlsx. Notes: [a] Results from previously published work following EUCAST breakpoints (14). [c] A _c_ denotes complement strand. [d] A transposase has inserted itself, splitting the *ermF* gene in two. Abbreviations: %ID; percent identity, %COV; coverage percentage, Chr; chromosome.

30

| Strain | Best Spades graph | Hybrid assembly after first nanopore sequencing run. No Canu correction | Hybrid assembly after first nanopore sequencing run with Canu correction | Hybrid assembly with data from supplementary runs with Canu correction | After manual finishing | Final assemblies contig statistics | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | No. | Length | Relative read depth |
| CCUG 4856T | | | | | | 1 | 5,205,133 | 1.00x |
| | | | | | | 2 | 36,560 | 7.42x |
| BFO17 | | | | | | 1 | 5,474,541 | 1.00x |
| | | | | | | 2 | 85,671 | 1.85x |
| | | | | | | 3 | 5,594 | 23.03x |
| BFO18 | | | | | | 1 | 5,302,644 | 1.00x |
| | | | | | | 2 | 7,221 | 25.98x |
| | | | | | | 3 | 4,137 | 50.80x |
| | | | | | | 4 | 2,782 | 59.91x |
| S01 | | | | | | 1 | 5,325,251 | 1.00x |
| | | | | | | 2 | 78,085 | 2.29x |
| | | | | | | 3 | 8,331 | 20.99x |
| | | | | | | 4 | 5,595 | 22.67x |
| BFO42 | | | | | | 1 | 5,141,257 | 1.00x |
| | | | | | | 2 | 8,306 | 40.06x |
| | | | | | | 3 | 5,594 | 40.15x |
| BFO67 | | | | | | 1 | 5,478,614 | 1.00x |
| | | | | | | 2 | 6,129 | 94.15x |
| BFO85 | | | | | | 1 | 5,504,076 | 1.00x |