

# Complete hybrid genome assembly of clinical multidrug-resistant *Bacteroides fragilis* isolates enables comprehensive identification of antimicrobial-resistance genes and plasmids

Thomas V. Sydenham<sup>1,2,3,\*</sup>, Søren Overballe-Petersen<sup>4</sup>, Henrik Hasman<sup>4</sup>, Hannah Wexler<sup>5</sup>, Michael Kemp<sup>1,2</sup> and Ulrik S. Justesen<sup>1,2</sup>

## Abstract

*Bacteroides fragilis* constitutes a significant part of the normal human gut microbiota and can also act as an opportunistic pathogen. Antimicrobial resistance (AMR) and the prevalence of AMR genes are increasing, and prediction of antimicrobial susceptibility based on sequence information could support targeted antimicrobial therapy in a clinical setting. Complete identification of insertion sequence (IS) elements carrying promoter sequences upstream of resistance genes is necessary for prediction of AMR. However, *de novo* assemblies from short reads alone are often fractured due to repeat regions and the presence of multiple copies of identical IS elements. Identification of plasmids in clinical isolates can aid in the surveillance of the dissemination of AMR, and comprehensive sequence databases support microbiome and metagenomic studies. We tested several short-read, hybrid and long-lead assembly pipelines by assembling the type strain *B. fragilis* CCUG4856<sup>T</sup> (=ATCC25285=NCTC9343) with Illumina short reads and long reads generated by Oxford Nanopore Technologies (ONT) MinION sequencing. Hybrid assembly with Unicycler, using quality filtered Illumina reads and Filtlong filtered and Canu-corrected ONT reads, produced the assembly of highest quality. This approach was then applied to six clinical multidrug-resistant *B. fragilis* isolates and, with minimal manual finishing of chromosomal assemblies of three isolates, complete, circular assemblies of all isolates were produced. Eleven circular, putative plasmids were identified in the six assemblies, of which only three corresponded to a known cultured *Bacteroides* plasmid. Complete IS elements could be identified upstream of AMR genes; however, there was not complete correlation between the absence of IS elements and antimicrobial susceptibility. As our knowledge on factors that increase expression of resistance genes in the absence of IS elements is limited, further research is needed prior to implementing AMR prediction for *B. fragilis* from whole-genome sequencing.

## DATA SUMMARY

Sequence read files [Oxford Nanopore Technologies (ONT) Fast5 files and Illumina FASTQ files], as well as the final genome assemblies, have been deposited in NCBI/ENA/DBJ under BioProject accession numbers PRJNA525024, PRJNA244942, PRJNA244943, PRJNA244944, PRJNA253771, PRJNA254401

and PRJNA254455. The FASTQ format of demultiplexed ONT reads trimmed of adapters and barcode sequences are available at <https://doi.org/10.5281/zenodo.2677927>. Genome assemblies from the assembly pipeline validation are available at <https://doi.org/10.5281/zenodo.2648546>. Genome assemblies corresponding to each stage of the process of the

Received 22 July 2019; Accepted 17 October 2019; Published 07 November 2019

**Author affiliations:** <sup>1</sup>Research Unit of Clinical Microbiology, Department of Clinical Research, University of Southern Denmark, Odense, Denmark; <sup>2</sup>Department of Clinical Microbiology, Odense University Hospital, Odense, Denmark; <sup>3</sup>Department of Clinical Microbiology, Lillebaelt Hospital, Vejle, Denmark; <sup>4</sup>Bacteria, Parasites and Fungi, Statens Serum Institut, Copenhagen, Denmark; <sup>5</sup>GLAVA Health Care System and David Geffen School of Medicine, UCLA (University of California, Los Angeles), Los Angeles, CA, USA.

\*Correspondence: Thomas V. Sydenham, Thomas.sydenham@rsyd.dk

**Keywords:** *Bacteroides fragilis*; antimicrobial resistance; genome sequencing; plasmid; Oxford Nanopore; hybrid assembly; insertion sequences.

**Abbreviations:** AMR, antimicrobial resistance; ANI, average nucleotide identity; CDS, coding sequence; % COV, per cent coverage; % ID, per cent identity; IQR, interquartile range; IS, insertion sequence; MDR, multidrug resistant; NCBI, National Center for Biotechnology Information; ONT, Oxford Nanopore Technologies; WGS, whole-genome sequencing.

Sequence files (MinION reads de-multiplexed with Deepbinner and basecalled with Albacore in Fast5 format, and Illumina MiSeq reads in FASTQ format) and final genome assemblies have been deposited in NCBI/ENA/DBJ under BioProject accession numbers PRJNA525024, PRJNA244942, PRJNA244943, PRJNA244944, PRJNA253771, PRJNA254401 and PRJNA254455.

**Data statement:** All supporting data, code and protocols have been provided within the article or through supplementary data files. Three supplementary figures and four supplementary tables are available with the online version of this article.

000312 © 2019 The Authors



This is an open-access article distributed under the terms of the Creative Commons Attribution License.

assembly are available at <https://doi.org/10.5281/zenodo.2661704>. Full commands and scripts used are available from GitHub (<https://github.com/thsyd/bfassembly>), as well as a static version (<https://doi.org/10.5281/zenodo.2683511>).

## INTRODUCTION

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

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

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

In 2014, we observed that identification of IS elements upstream of known AMR genes in *B. fragilis* was hampered in short-read *de novo* assemblies even though the genes could be identified [14]. This occurred because contigs were often terminated close to the start of the resistance genes, presumably due to the proliferation of multiple copies of the same IS elements throughout the *B. fragilis* genomes. Genome assemblies from short-read sequencing technologies alone most often result in fragmented assemblies,

## Impact Statement

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

because of repetitive regions and genome elements with multiple occurrences in the chromosomes and plasmids [15, 16]. Therefore, we could not predict AMR phenotypes in *B. fragilis* using only short reads for WGS, since IS element identification is a prerequisite for correct genotype–phenotype associations. Long-read sequencing technologies are increasingly being utilized to increase the contiguity of bacterial genome assemblies, and often result in complete, closed chromosomes and plasmids [17–20]. This provides possibilities for comprehensive identification of IS elements,

insights into genome structures, and characterization of other mobilizable elements and associated genes. Complete identification and characterization of plasmids in sequenced isolates would allow for improved analysis of the plasmid-mediated spread of AMR.

Bioinformatic analysis of WGS data depends heavily on high-quality reference databases. Anaerobes make up most of the bacterial human commensal microbiota, but are most likely underrepresented in public databases of whole genomes from cultured isolates. The National Center for Biotechnology Information (NCBI) genome database (accessed 31/03/2019) contains genome sequences of 191 411 bacteria, of which 13 483 are marked as complete assemblies. Only seven of these are *B. fragilis* [21–27]. In comparison, there are 776 assemblies of *Escherichia coli* marked as complete and 398 of *Staphylococcus aureus*. Improving the representation of complete assemblies of *B. fragilis* in the public genome databases will support the development of AMR prediction from WGS, as well as microbiome and metagenomic analysis projects.

The aims of this study were to select an optimal assembly software pipeline for complete, circular assembly of *B. fragilis* and demonstrate the utility of complete assembly for both plasmid identification and comprehensive detection of genes and IS elements associated with AMR. We assembled the *B. fragilis* CCUG4856<sup>T</sup> (=ATCC25285=NCTC9343) reference strain utilizing long reads generated with the MinION sequencer from Oxford Nanopore Technologies (ONT) and high-quality Illumina short reads, and selected the best assembly pipeline by comparing assemblies to the Sanger-sequenced reference NCTC9343 (RefSeq accession no. GCF\_000025985.1). The best assembly pipeline was then applied to six clinical MDR *B. fragilis* isolates from our 2014 study [14].

## METHODS

### Culture conditions and DNA extraction

*B. fragilis* CCUG4856<sup>T</sup> and the six strains described in our previous study were included [14, 21]. Strains were stored at –80°C in beef extract broth with 10 % (v/v) glycerol (SSI Diagnostica), and cultured on solid chocolate agar with added vitamin K and cysteine (SSI Diagnostica) for 48 h in an anaerobic atmosphere at 35 °C. Ten microlitres of culture was transferred to 14 ml saccharose serum broth (SSI Diagnostica) and incubated for 18 h under the same conditions. DNA was then extracted using the Genomic-Tip G/500 kit (Qiagen) following the manufacturers protocol for Gram-negative bacteria and eluted into 5 mM Tris pH 7.5, 0.5 mM EDTA buffer. Quality control was performed by measuring fragment length on a TapeStation 2500 (Genomic DNA ScreenTape; Agilent), purity on a NanoDrop instrument (ThermoFisher Scientific) and concentration on a Qubit instrument (dsDNA BR kit; Invitrogen). The eluted DNA was then stored at –20 °C.

### Illumina library preparation, sequencing and quality control

The strains had previously been sequenced and assembled using Illumina short reads for our previous study [14], but to minimize biological disparities we opted to re-sequence with Illumina using the same DNA extraction prepared for long-read sequencing. Paired-end libraries were generated using the Nextera XT DNA sample preparation kit (Illumina), according to the manufacturer's protocol. DNA was sequenced on a MiSeq sequencer (Illumina) with 150 bp reads for a theoretical read depth of 100×. Read quality metrics were evaluated using FASTQC (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and FASTP v0.19.6 [28]. Filterbytile from the BMap package (<http://sourceforge.net/projects/bmap/>) was used with default parameters for removing low-quality reads based on positional information on the sequencing flowcell and Trim Galore ([http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)), with settings --qual 20 and --length 126, provided additional adapter and quality trimming. FASTQ files were then randomly down-sampled to <100× crude read depth using an estimated genome size of 5.3 Mb, as higher read depths tend to reduce assembly quality [29].

### Nanopore library preparation and MinION sequencing

Sequencing libraries were prepared using the Rapid Barcoding kit (SQK-RPB004; ONT) following the manufacturers protocol (version RPB\_9059\_v1\_revC\_08Mar2018) with SPRI bead clean up (AMPure XT beads; Beckman Coulter) as described. Sequencing was performed as multiplex runs on a MinION connected to a Windows PC with MinKnow v1.15.1 using FLO-MIN106 R9.4 flowcells. Raw Fast5 files were transferred to the Computerome high-performance cluster (<https://www.computerome.dk/>) for analysis. Four sequencing runs were performed, as the first two runs did not provide enough data for complete assembly of all isolates (see Results).

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

The raw Fast5 files were demultiplexed with Deepbinner v0.2.0 and base-called using Albacore v2.3.3, retaining only those barcodes Deepbinner and Albacore agreed upon for minimal barcode misclassification [30]. Porechop v0.2.4 (<https://github.com/rrwick/Porechop>) with --check\_reads 100 and --discard\_middle options was used for adapter and barcode trimming, and read statistics were collected using NanoPlot [31]. Filtlong v0.2.0 (<https://github.com/rrwick/Filtlong>), with parameters --min\_length 100 --keep\_percent 90 --target\_bases 500000000, was used to filter the long reads by either removing the worst 10 % or by retaining 500 Mbs in total, whichever option resulted in fewer reads.

### Assembly validation

To select and validate the optimum assembly pipeline *B. fragilis* CCUG4856<sup>T</sup> was assembled using a variety of well-known

**Table 1.** Genome assemblers and polishing tools tested

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

assemblers and polishing tools (Table 1). Each assembler was run with the Filtrong filtered reads as input or the filtered reads corrected with Canu 1.8 (with `--genomSize=5.4 m` and `corMinCoverage=0` or `coroutCoverage=999`). Canu was also tested with the unfiltered reads as input. Hybrid assemblers used the filtered long reads and the filtered, trimmed and down-sampled Illumina reads. Unicycler includes polishing with Racon and Pilon. For assemblers other than Unicycler, Racon polishing with ONT reads was run for one or two rounds, and Pilon was run until no changes were made or for a maximum of six rounds. Racon polishing with Illumina reads was run for one round.

The original Sanger-sequenced *B. fragilis* NCTC9343 (=CCUG4856<sup>T</sup>) [21] downloaded from NCBI RefSeq (accession no. GCF\_000025985.1) was used as a reference sequence for the assembly comparisons and Quast v5.0.2 was used for assembly summary statistics, indel count and K-mer-based completion [32]. BUSCO v3.0.2b (with parameters: `--lineage_path <path to the bacteroidetes_odb9dataset> --mode genome --force`), CheckM v1.0.12 (with default parameters) and Prokka v1.13.3 (with parameters: `--compliant`) were used to assess gene content [33–35]. Average nucleotide identity (ANI) was calculated using software available at <https://github.com/chjp/ANI/blob/master/ANI.pl> and ALE v0.9, which uses a likelihood-based approach to assess the quality of different assemblies based on alignment of Illumina reads, was also used to score the assemblies using default parameters [36, 37]. Ranking of assemblies was based on number of contigs, number of circular contigs, closeness to total length compared to the reference genome, number of local misassemblies, number of mismatches per 100 kb, number of indels

per 100 kb, ANI, CheckM and BUSCO scores, and the total ALE score (a higher score is better). Please see <https://github.com/thysyd/bfassembly> for full bioinformatics methods.

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

The assembly strategy deemed to produce the highest-quality genome for CCUG4856<sup>T</sup> was chosen for initial assembly of the six MDR *B. fragilis* isolates. Manual finishing of incomplete assemblies was performed using Bandage for visualization of assembly graphs and blastn searches [38]. Minimap2 and BWA MEM were used to map reads to the assemblies for coverage graphs [39, 40]. Long-read assembly with Flye was compared to the Unicycler assembly, and used to guide and validate the manual finishing results. Circlator's *fixstart* task was used to fix the start position of the manually finished genomes to be at the *dnaA* gene [41].

The assembled genomes were submitted to NCBI GenBank and annotated with PGAP [42]. ABRicate v0.8.10 (<https://github.com/tseemann/ABRicate>) (with options `--minid 40 --mincov 25`) was used to screen for AMR genes with the ResFinder (database date 19/08/2018), NCBI Bacterial Antimicrobial Resistance Reference Gene Database (database date 19/09/2018) and CARD (v2.0.3) databases, supplemented with nucleotide sequences for the multidrug efflux-pump genes *bexA* (GenBank accession no. AB067769.1: 3564...4895) and *bexB* (GenBank accession no. AY375536.1: 4599...5963) [43, 44]. IS elements were identified 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

The PLSDb web server (<https://ccb-microbe.cs.uni-saarland.de/plsdb/>) (data v. 2019\_03\_05) contains bacterial plasmid sequences retrieved from the NCBI, and was used for screening and identifying putative plasmids sequences [46]. Only hits to accessions from cultured organisms were included. Putative plasmids not identified using PLSDb were evaluated by the read depth relative to the chromosome (higher relative read depth indicates plasmid sequence) and Pfam families covering known plasmid replication domains from table 1 in the 2014 reference by Jørgensen and colleagues [47] were downloaded from the Pfam database (Pfam 32.0; <https://pfam.xfam.org/>) and used for screening putative plasmids with ABRicate.

## RESULTS

### Sequencing data quality

For Illumina data, a median of 3 465 082 reads [interquartile range (IQR): 3 177 493–5 001 077] were generated for each isolate (Table S1, available with the online version of this article). After filtering, adapter removal and down sampling, a median of 449 022 741 bases (IQR: 433 517 549–530, 57 210) was available per isolate with 87–96 % Q30 bases corresponding to calculated read depths of 75–103 %. The mol%



G+C content of the reads for each isolate (median 42.9 mol%, range 42.6–43.3 mol%) were very consistent and within the expected range for the genus *Bacteroides* (40–48 mol%) [48].

Isolates were sequenced in runs multiplexed with other isolates not included in this study. Based on initial test assemblies using Unicycler without filtering or Canu correction (not shown), it was concluded that data from the first ONT sequencing runs should be supplemented by additional runs to increase the chance of complete assembly of all isolates. Concatenating reads from runs, a median of 75 598 reads (IQR: 50 210–112 065) with a median length of 2938–4393 bases were generated for each isolate (Table S1). Filtering with FilTlong and correction with Canu resulted in a median of 8515 reads (IQR: 6226–10 370) with median lengths of 6181–38 588 bases for each isolate as input for the assemblies.

### Selecting the optimal assembly pipeline

A total of 141 assemblies of *B. fragilis* CCUG4856<sup>T</sup> was generated using the various assemblers and polishing steps (Table S2). Compared to the reference genome, Unicycler assemblies were of the highest quality (Table 2). Unicycler, with any of the read input options, produced two circular contigs of the expected lengths, and the differences between the various Unicycler assemblies were minimal (Table 3). Assemblies with Canu-corrected reads showed slightly higher genome fractions and ANIs to the reference and fewer mismatches and indels, when compared to Unicycler alone. Unicycler assemblies corrected with Racon using Illumina reads worsened slightly overall with 0.04–0.19 more indels and 0.14–0.25 more mismatches per 100 kbp. Based on this initial evaluation, the assembly pipeline using Canu-corrected reads with default options was chosen (assembly 'OFCS' in Table 3). This would reduce the number of long reads, compared to Canu correction with corMinCoverage=0 or coroutCoverage=999, and thereby lead to a faster run-time for Unicycler.

The hybrid Unicycler assembly of CCUG4856<sup>T</sup> with standard Canu-corrected ONT reads consists of two circular contigs of 5 205 133 and 36 560 bp in length. The plasmid is the same length as plasmid pBF9343 from the reference assembly GCF\_000025985.1 and the chromosome is seven bases shorter. Alignments of the Sanger-sequenced assembly GCF\_000025985.1 with the hybrid Unicycler assembly show an 88045 bp inversion in the hybrid assembly compared to the Sanger assembly (Fig. 1). This inversion is present in all the best assemblies, including assemblies derived from solely ONT sequences or Illumina sequences (Fig. S1), as well as two additional assemblies of NCTC9343/ATCC25285 from PacBio and Illumina sequences downloaded from NCBI RefSeq (Fig. S2).

### Complete assembly of six MDR isolates

Unicycler, using filtered and trimmed Illumina reads and the FilTlong filtered and Canu-corrected ONT reads from the first sequencing runs, generated complete, continuous, circular assemblies for two of the six isolates (BFO18 and BFO67) (Fig. 2). For the assemblies that were not complete with

sequencing data from the first MinION runs, increasing the amount of ONT data resulted in fewer contigs overall, except for BFO67, where the additional data from the second sequencing run led to a fragmented assembly and manual finishing was necessary. Performing assembly of isolate S01 without Canu correction of the ONT reads from the first sequencing resulted in a closed chromosome and performing Canu correction of reads resulted in a fragmentation of the chromosome. This was ameliorated by including more ONT data. By manual finishing using read mapping and additional assembly with Flye, the remaining three assemblies were circularized. Chromosomes varied in length from 5 141 257 to 5 504 076 bp. Alignment of ONT and Illumina reads to the chromosome assemblies showed even coverage for both sequencing technologies (Fig. S3). For BFO85, a >100 % relative read depth increase was observed at approximately 25–38 kb. This could represent a 12 kb repeat region that was not resolved in the assembly. Seven (47 %) of the fifteen PGAP annotated coding sequences (CDSs) in the 13 kb region were annotated as hypothetical proteins. None of the annotated CDSs represented mobilizable proteins.

### Eleven putative plasmid sequences were identified

A total of 11 putative circular plasmids were identified in the six *B. fragilis* isolates (Table 4). Zero to three putative plasmids were identified per isolate with lengths varying from 2782 to 85 671 bp.

The PLSDb database contains NCBI RefSeq plasmid sequences marked as complete. Three of the eleven putative plasmid sequences were found to match (ID >98 %) a sequence in PLSDb (Table 4). These three all matched the cryptic plasmid pBFP35 [49]. The NCBI nucleotide database was queried using BLASTN with the remaining unidentified putative plasmid sequences [50]. BFO18 putative plasmid sequence pBFO18\_1 (7221 bp) resembles plasmid pIP421, a 7.2 kb plasmid with metronidazole-resistance gene *nimD* and IS1169. Partial sequences in NCBI GenBank spanning the *nimD* gene, IS element and *repA* (GenBank accession numbers Y10480.1 and X86702.1) showed 99 % ID (per cent identity) to their alignment to pBFO18\_1 (not shown) [51, 52]. Strain S01 putative plasmid sequence pBFS01\_2 (8331 bp) showed 99.87 % ID to the 1486 bp partial sequence of *B. fragilis* plasmid pBF388c (GenBank accession number AM042593.1), an 8.3 kb conjugative plasmid harbouring *nimE* and ISBf6 [53].

None of the three putative plasmid sequences of strain BFO18 could be identified using the PLSDb, but querying the NCBI nucleotide database using BLASTN revealed hits for all three. The hits corresponded to circularized sequences [% ID, 99.56–99.96; per cent coverage (% COV), 100] assembled from mobilome metagenomic sequencing of the uncultured caecum content from a rat trapped at Bispebjerg Hospital in Copenhagen, Denmark (a 2 h drive from Odense University Hospital where BFO18 was isolated from a patient's blood culture) (Table S3) [47, 54]. BLASTN searches of the remaining unidentified putative plasmids from the other strains did not reveal complete hits.

**Table 2.** Selected quality indicators for the best genome assembly of *B. fragilis* CCUG4856<sup>T</sup> per assembly pipeline

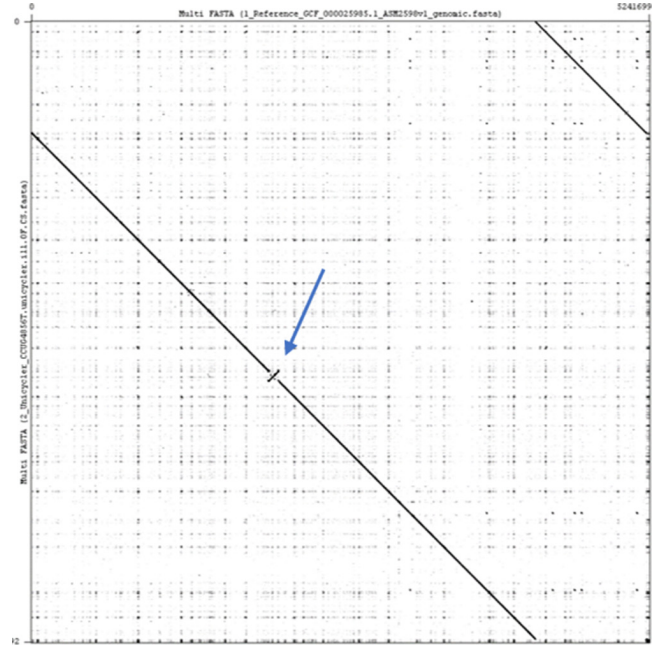
RefSeq accession GCF\_000025985.1 was used as a reference. CM, Canu-corrected with option corMinCoverage=0; CO, Canu-corrected with option corOutCoverage=999; CS, Canu-corrected standard settings; OF, ONT reads filtered with FilTlong; P[n], Pilon polishing with Illumina reads, [n] rounds; RI, Racon polishing with Illumina reads, [n] rounds; RO2, two rounds of Racon polishing with ONT reads. Full results are available in Table S2.

Assembly	No. of contigs	Largest contig	Total length	Mis-assemblies	Genome fraction (%)	Mismatches per 100 kbp	Indels per 100 kbp	ANI	CheckM completeness	Busco score: 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	100.000	99.26	442/0/1	4439	19	73	-17071758.95
Skesa	46	553 341	5 201 945	3	99.237	0.23	0.15	99.998	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	99.987	99.26	440/2/1	4407	3	56	-19676529.39
Canu.OFCO.RO2.RI.PI3	2	5 247 938	5 350 432	8	99.972	4.94	15.9	99.975	99.26	442/0/1	4634	19	73	-19283611.73
Flye.OFCO.PI5.RI	5	2 282 650	5 269 269	4	99.917	1.07	6.24	99.978	99.26	441/1/1	4476	19	73	-18222322.23
Miniasm.OFCM.RO2.PI5	3	5 204 445	5 277 434	2	99.972	5.21	17.75	99.969	98.88	442/0/1	4607	19	73	-17789234.97
Wtdbg2.OFCO.RO2.PI6.RI	3	5 192 352	5 234 448	7	99.723	3.23	3.04	99.981	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	99.986	99.26	440/2/1	4441	19	73	-18535980.68
Unicycler.OFCS	2	5 205 133	5 241 693	2	99.972	0.84	0.48	100.000	99.26	442/0/1	4435	19	73	-17200232.52

**Table 3.** Hybrid Unicycler assemblies of *B. fragilis* CCUG4856<sup>T</sup>

RefSeq accession GCF\_000025985.1 was used as a reference. CM, Canu-corrected with option corMinCoverage=0; CO, Canu-corrected with option corOutCoverage=999; CS, Canu-corrected standard settings; OF, ONT reads filtered with Filtlong; RI, Racon polishing with Illumina reads. Unicycler performs assembly polishing with Racon (ONT reads) and Pilon. Full results are available in Table S2.

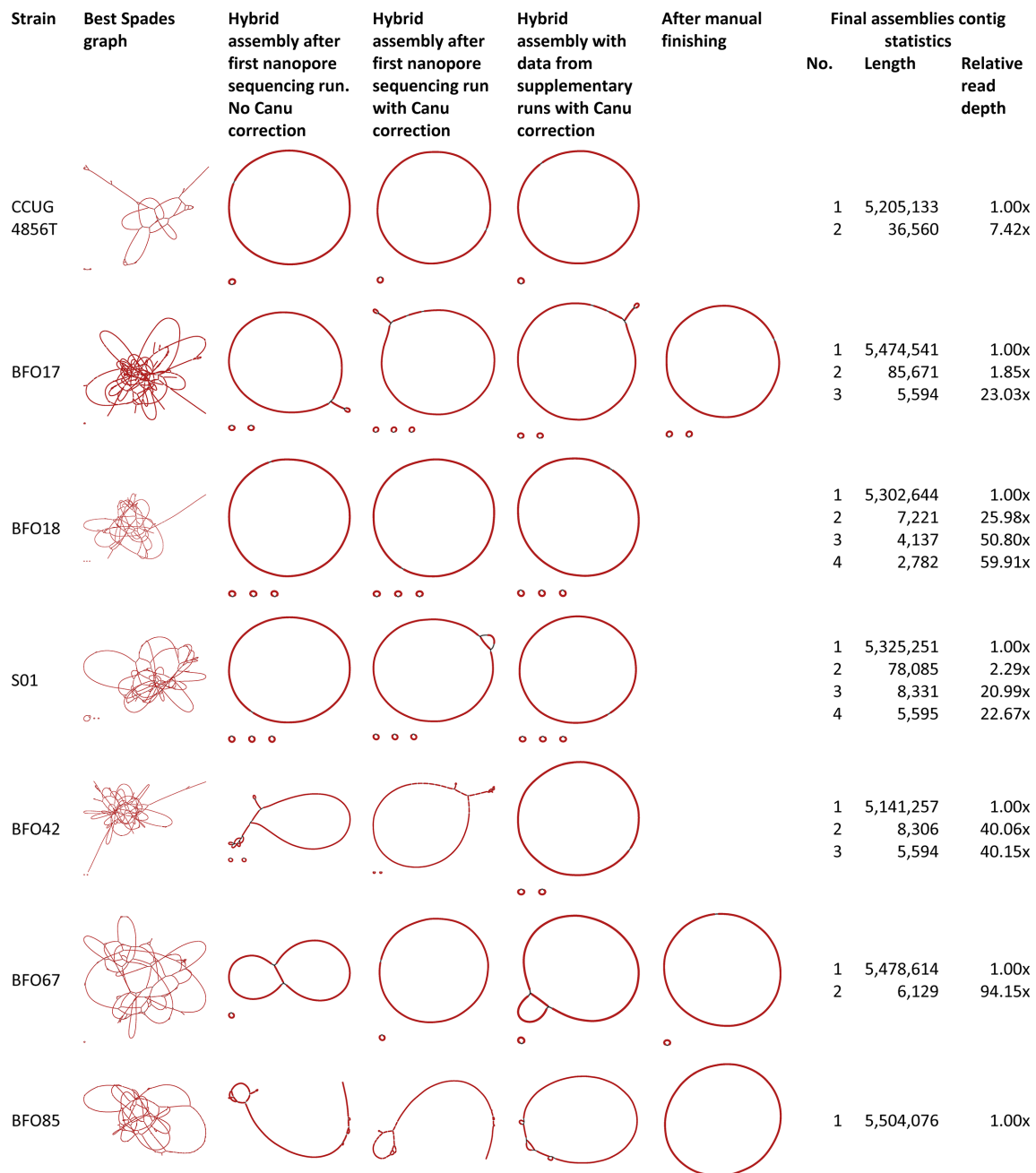
Assembly	Total length (bp)	Largest contig (bp)	Local mis-assemblies	Genome fraction (%)	Mismatches per 100 kbp	Indels per 100 kbp	K-mer-based compl. (%)	K-mer-based misjoins	ANI	Prokka CDSs	Prokka genes	Total ALE score
GCF_000025985.1	5 241 700	5 205 140	0	100.000	0	0	100.00	0	100.000	4346	4439	-17071758.95
OF	5 241 602	5 205 042	3	99.970	1.11	0.65	99.96	0	99.999	4343	4436	-17245134.52
OF.RI	5 241 606	5 205 046	3	99.970	1.09	0.67	99.96	3	99.999	4345	4438	-17247815.86
OF.CS	5 241 693	5 205 133	2	99.972	0.84	0.48	99.97	1	100.000	4342	4435	-17200232.52
OF.CS.RI	5 241 698	5 205 138	2	99.972	0.88	0.52	99.96	1	100.000	4346	4439	-17206271.66
OF.CM	5 241 691	5 205 131	2	99.972	0.88	0.5	99.96	1	100.000	4343	4436	-17201292.44
OF.CM.RI	5 241 696	5 205 136	2	99.972	0.95	0.55	99.97	1	100.000	4343	4436	-17193184.79
OF.CO	5 241 693	5,205,133	2	99.972	0.84	0.48	99.97	1	100.000	4342	4435	-17200232.52
OF.CO.RI	5 241 698	5 205 138	2	99.972	0.88	0.52	99.96	1	100.000	4346	4439	-17206271.66



**Fig. 1.** Dot plot matrix of the alignment of the reference assembly and the hybrid Unicycler assembly using Gepard v1.40 [81]. The *B. fragilis* NCTC9343 (RefSeq accession number GCF\_000025985.1) reference assembly derived from Sanger sequencing is on the x-axis and the hybrid Unicycler assembly on the y-axis. On this otherwise near-perfect alignment with high similarity, an 88 045 bp inversion with 100 % ID is observed at nucleotide positions 2 941 962...3 030 006 on the Unicycler assembly (2 005 742...2 093 786 on the reference sequence) (indicated by the blue arrow).

Using ABRicate with the plasmid replication domains collected from the Pfam database, all putative plasmids, except pBF017\_1 and pBFS01\_1, were found to have recognized replicon domains (Table 4). The circular structures of the two sequences lacking a predicted replication domain were confirmed manually by visually inspecting BLASTN mapping of ONT sequences longer than 10 kbp to the assembled plasmid sequences with CLC Genomics Workbench 10 (Qiagen). A total of 11 and 22 ONT reads spanned the complete lengths of pBF017\_1 and pBFS01\_1, respectively, and contained no other elements. pBF017\_1 and pBFS01\_1 demonstrate a degree of similarity of close to 100%, except for an approximate total of 7500 bp transposase and prophage sequences in pBF017\_1 (Fig. 3). No alignment to chromosomal sequences of any of the included *B. fragilis* isolates was observed using progressiveMauve (not shown) [55].

The G+C contents of pBF017\_1 and pBFS01\_1 are 36.78 and 36.04 mol%, respectively. These lie within the range for the genus *Bacteroides* but differ from the expected value for *B. fragilis* (43 mol%), which could indicate that the putative plasmids do not originate from *B. fragilis* [56]. After supplementing the PGAP annotations with RAST annotation [57], 63 % (pBF017\_1) and 59 % (pBFS01\_1) of CDSs remained annotated as hypothetical or as domain of unknown function. Of the annotated CDSs, the majority were associated



**Fig. 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 circularized 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 Table S1.

with mobilizable features, plasmids and phages such as *parA* and *parB*, DNA partitioning proteins, conjugative transposon proteins, transposases, DNA binding motif domain containing proteins, and reverse transcriptase protein. The results above support the assembly data suggesting these two sequences are in fact plasmids.

### Detection of AMR genes and IS 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 homologues to known AMR genes and IS elements adjunct to the AMR genes, were



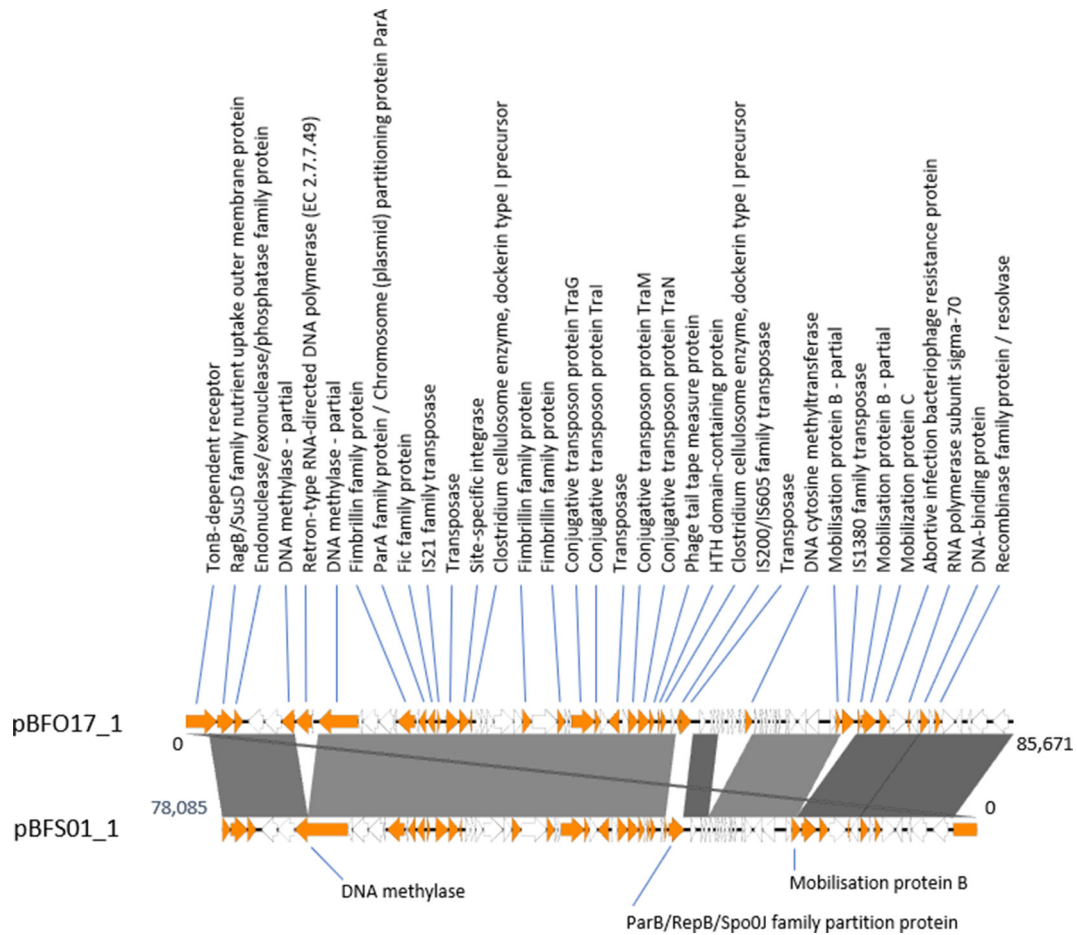
**Table 4.** Putative plasmid sequences of the complete *B. fragilis* assemblies

Putative plasmid sequences from the hybrid assemblies of *B. fragilis* CCUG4856<sup>T</sup> and the six MDR *B. fragilis* isolates were screened using the PLSDB. The best hit to plasmids from cultured isolates is shown. Only three putative plasmids from the MDR *B. fragilis* isolate assemblies could be identified with confident % ID. For most sequences, plasmid replication family proteins were identified in the putative plasmids using ABRicate with a database of sequences downloaded from the Pfam database, strengthening the interpretation that the circularized putative plasmid sequences do in fact represent plasmids harboured by the isolates.

Strain	Sequence	Length (bp)	Relative read depth	mol% G+C	PLSDB results			Plasmid replicon family (% COV, % ID)
					Best hit accession no.	Plasmid hit name	% ID	
CCUG4856 <sup>T</sup>	Chr	5205133	1.00×	43.19	-	-	-	-
	pBF9343	36560	7.42×	32.19	NC_006873.1	pBF9343	100	Rep_3 (100/100)
BFO17	Chr	5474541	1.00×	43.51	-	-	-	-
	pBFO17_1	85671	1.85×	36.78	NC_006873.1	pBF9343	80.7	None
BFO18	pBFO17_2	5594	23.03×	39.65	NC_011073.1	pBFP35	99.9	Rep_1 (100/100)
	Chr	5302644	1.00×	43.34	-	-	-	-
S01	pBFO18_1	7221	25.98×	42.32	NC_015168.1	pBACSA02	85.6	Rep_3 (99.69/99.69)
	pBFO18_2	4137	50.80×	45.40	NC_019534.1	pBFUK1	92.2	Rep_3 (100.00/98.24)
BFO42	pBFO18_3	2782	59.91×	41.45	NC_005026.1	pBI143	94.6	Repl. (89.66/49.22)*
	Chr	5325251	1.00×	43.57	-	-	-	-
BFO67	pBF801_1	78085	2.29×	36.04	NC_006873.1	pBF9343	80.7	None
	pBF801_2	8331	20.99×	41.17	NC_015166.1	pBACSA03	95.6	Rep_3 (100.00/97.85)
BFO85	pBF801_3	5595	22.67×	39.62	NC_011073.1	pBFP35	99.9	Rep_1 (100.00/99.48)
	Chr	5141257	1.00×	43.35	-	-	-	-
BFO85	pBFO32_1	8306	40.06×	43.34	KJ830768.1	pBF69566b	96.0	RHH_1 (92.94/64.63) Rep_3 (93.64/68.31)
	pBFO32_2	5594	40.15×	39.63	NC_011073.1	pBFP35	99.9	Rep_1 (100.00/99.48)
BFO85	Chr	5478614	1.00×	43.85	-	-	-	-
	pBFO67_1	6129	94.15×	41.67	NC_011073.1	pBFP35	76.9	Rep_3 (100.00/99.69)
BFO85	Chr	5504076	1.00×	43.60	-	-	-	-

Chr, Chromosome.

\*Annotated as RepA protein in the PGAP annotation.



**Fig. 3.** Linear representation of an alignment of putative circular plasmid sequences pBF017\_1 and pBFS01\_1 (reverse complement for better visualization) using EasyFig [82]. EasyFig uses BLAST to identify sequences of similarity. Sequence similarities of >98 % are indicated by full colouring, a darker colour indicates a higher % ID. Products of annotated CDSs are shown. CDSs annotated as hypothetical or domain of unknown function are coloured white. The two sequences show a very high degree of similarity. pBF017\_1 is 7586 bp longer than pBFS01\_1. This is mainly due to the insertion of a reverse transcriptase (pBF017\_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 IS1380 family-like transposase (67933...69237). The regions pBF017\_1 50711...52501 and pBFS01\_1 32248...30304 are not similar. Possibly, the insertion of two transposases in pBF017\_1 have excised most of the ParB-family DNA partitioning protein in the corresponding sequence range in pBFS01\_1.

detected (Table 5). Of note, isolate BFO17 contains two homologues of the metronidazole-resistance gene *nimJ* (with a 100 % consensus) and two isolates, S01 and BFO85, harbour two homologues of the tetracycline-resistance gene *tetQ*. Homologues 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* were observed for several isolates, but with low % ID and % COV, and possibly represent identification of conserved domains, but not *ugd* homologues. 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 *B. fragilis* promoters 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. Minimum inhibitory concentrations (MICs) for meropenem

**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 homologues for genes encoding PmrE, which is involved in polymyxin resistance in Gram-negative bacteria. Full ABRicate results with nucleotide positions and information on the IS elements are available in Table S4.

Antimicrobial susceptibility*				AMR genes and IS elements						
Strain	Antimi-crobia	Etest MIC (mg l <sup>-1</sup> )	Result	Gene	Upstream IS element	Sequence†	% ID	% COV	Associated resistance to drug class	
BFO17	MEM	>32	R	<i>cfiA11</i>	IS614B	Chr	100.00	99.20	Carbapenem	
	IPM	>32	R							
	MTZ	>32	R	<i>nimJ</i>	IS614B	Chr	99.40	100.00	Nitroimidazole	
				<i>nimJ</i>	IS614B	Chr <sub>c</sub>	99.40	100.00	Nitroimidazole	
	CLI	0.094	S							
	TZP	>256	R							
				<i>tetQ</i>		Chr	99.34	99.34	Tetracycline	
				<i>cfxA4</i>		Chr	85.71	100.00	Cephamicin	
				<i>bexB</i>		Chr <sub>c</sub>	91.21	100.00	Fluoroquinolone	
				<i>bexA</i>		Chr	73.77	99.02	Fluoroquinolone	
BFO18	MEM	>32	R	<i>cfiA2_1</i>	ISBf12	Chr	100.00	100.00	Carbapenem	
	IPM	16	R				100.00	100.00		
	MTZ	16	R	<i>nimD</i>	IS1169	S	99.19	100.00	Nitroimidazole	
	CLI	6	R	<i>ermF</i> ‡	IS4351	Chr <sub>c</sub>	99.83	72.03	Clindamycin	
					ISBthe1‡	Chr <sub>c</sub>	70.97	97.19		
				<i>erm(F)</i> §		Chr <sub>c</sub>	99.58	29.71		
				<i>lnu(AN2)</i>		Chr <sub>c</sub>	100.00	100.00	Clindamycin	
	TZP	>256	R							
				<i>ugd</i>		Chr	65.69	53.04	Polymyxin	
				<i>bexA</i>		Chr <sub>c</sub>	73.60	99.02	Fluoroquinolone	
				<i>bexB</i>		Chr	91.14	100.00	Fluoroquinolone	
				<i>tet(Q)</i>		Chr <sub>c</sub>	99.79	100.00	Tetracycline	
				<i>mef(En2)</i>		Chr <sub>c</sub>	99.83	100.00	Macrolides	
S01	MEM	>32	R	<i>cfiA13_1</i>	IS1187	Chr	99.20	100.00	Carbapenem	
	IPM	16	R							
	MTZ	64	R	<i>nimE</i>	ISBf6	pBFS01_2	100.00	100.00	Nitroimidazole	
	CLI	>32	R	<i>erm(F)</i>	IS1187	Chr	99.50	100.00	Clindamycin	
	TZP	6	S							
				<i>tetQ</i>		Chr	90.02	99.95	Tetracycline	
				<i>tet(Q)</i>		Chr <sub>c</sub>	99.84	100.00	Tetracycline	
				<i>bexB</i>		Chr <sub>c</sub>	91.06	100.00	Fluoroquinolone	
				<i>bexA</i>		Chr	74.03	98.80	Fluoroquinolone	
BFO42	MEM	0.094	S							

Continued

Table 5. Continued

Antimicrobial susceptibility*				AMR genes and IS elements					
Strain	Antimicrobials	Etest MIC (mg l <sup>-1</sup> )	Result	Gene	Upstream IS element	Sequence†	% ID	% COV	Associated resistance to drug class
BFO67	IPM	0.25	S						
	MTZ	8	R	<i>nimA</i>	IS <i>Bf13</i>	pBFO32_1	98.64	96.61	Nitroimidazole
	CLI	>256	R	<i>erm(F)</i>	IS <i>613</i>	Chr	99.50	100.00	Clindamycin
				<i>lnu(AN2)</i>		Chr	100.00	100.00	Clindamycin
	TZP	0.38	S	<i>ugd</i>		Chr	70.38	31.45	Polymyxin
				<i>cepA-49</i>		Chr <sub>c</sub>	100.00	100.00	Cephalosporin
				<i>mef(En2)</i>		Chr	99.83	100.00	Macrolide
				<i>ugd</i>		Chr	71.15	31.11	Polymyxin
				<i>tetQ</i>		Chr <sub>c</sub>	100.00	100.00	Tetracycline
				<i>bexB</i>		Chr	99.12	100.00	Fluoroquinolone
				<i>ere(D)</i>		Chr	96.66	100.00	Erythromycin
				<i>aadS</i>		Chr <sub>c</sub>	99.88	100.00	Aminoglycoside
				<i>OXA-347</i>	IS <i>613</i>	Chr <sub>c</sub>	100.00	100.00	Penicillin, cephalosporin
	<i>bexA</i>		Chr	75.09	99.62	Fluoroquinolone			
	MEM	8	R	<i>cfiA13_1</i>	None	Chr	100.00	100.00	Carbapenem
	IPM	0.5	S						
	MTZ	0.19	S						
	CLI	0.38	S						
	TZP	2	S	<i>cfxA2</i>	IS <i>Bf11</i>	Chr	99.69	100.00	Cephameycin
<i>mef(En2)</i>					Chr	99.75	100.00	Macrolide	
<i>lnu(AN2)</i>					Chr	100.00	100.00	Clindamycin	
<i>ugd</i>					Chr <sub>c</sub>	66.76	56.30	Polymyxin	
<i>tet(Q)</i>					Chr	100.00	100.00	Tetracycline	
<i>bexB</i>					Chr <sub>c</sub>	90.92	100.00	Fluoroquinolone	
<i>bexA</i>					Chr	73.90	99.02	Fluoroquinolone	
BFO65	MEM	32	R	<i>cfiA2_1</i>	None	Chr	100.00	100.00	Carbapenem
	IPM	1	S						
	MTZ	0.25	S						
	CLI	>256	R	<i>ermB</i>		Chr <sub>c</sub>	99.19	98.66	Clindamycin
	TZP	2	S	<i>ugd</i>		Chr	69.84	31.45	Polymyxin
				<i>tetQ</i>		Chr <sub>c</sub>	90.02	99.95	Tetracycline
<i>aadE</i>					Chr <sub>c</sub>	100.00	100.00	Aminoglycoside	
<i>aad9</i>		Chr <sub>c</sub>	100.00	100.00	Aminoglycoside				

Continued



Table 5. Continued

Antimicrobial susceptibility*				AMR genes and IS elements					
Strain	Antimi-crobial	Etest MIC (mg I <sup>-1</sup> )	Result	Gene	Upstream IS element	Sequence†	% ID	% COV	Associated resistance to drug class
				<i>bexB</i>		Chr <sub>c</sub>	90.92	100.00	Fluoroquinolone
				<i>bexA</i>		Chr	73.53	99.02	Fluoroquinolone
				<i>cfxA2</i>	IS614	Chr <sub>c</sub>	100.00	100.00	Cephameycin
				<i>tet(Q)</i>		Chr <sub>c</sub>	99.84	100.00	Tetracycline

Chr, Chromosome; MEM, meropenem; IPM, imipenem; MTZ, metronidazole; CLI, clindamycin; TZP, piperacillin/tazobactam.

\*Results from previously published work following EUCAST (European Committee on Antimicrobial Susceptibility Testing) breakpoints [14].

†A subscript letter C denotes the complement strand.

‡A transposase has inserted itself, splitting the *ermF* gene in two.

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 but for one isolate (BFO85) an upstream IS element was not found. 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 I<sup>-1</sup>) 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* CCUG4295<sup>T</sup> with high similarity to the reference assembly of the original Sanger-sequenced reference assembly. An 88 kb 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 strains stored at the National Collection of Type Cultures (NCTC) and the Culture Collection University of Gothenburg (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 88 kb sequence.

### Complete genome assembly of three of the six MDR isolates required manual finishing

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

ONT read depth did not serve as an indicator of whether the Unicycler assemblies would result in closed chromosomal contigs in this study. Final ONT read depth, prior to Filtlong

filtering and Canu correction, ranged from 23–371×, but a high read depth alone was not an indicator of closed contigs. The three assemblies BFO17, BFO67 and BFO85 required manual finishing to complete the assemblies, and had ONT raw read depths of 99–137×. After Filtlong filtering and Canu correction, the median read lengths were 21 932–29 893 bases and read length N50 was 25 765–34 815 bases for the three isolates (Table S1). Canu correction improved the Unicycler assembly of *B. fragilis* CCUG4856<sup>T</sup> by nearly all parameters. But whilst Canu correction of the data from the first sequencing run resulted in the complete assembly of BFO67, the assembly of S01 worsened slightly. Increasing the amount of ONT data for BFO67 fragmented the complete chromosome. However, increasing the ONT read depth did decrease the number of contigs per isolate in our study overall.

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

We chose to benchmark a selection of widely used genome assemblers for short-read, long-read and hybrid bacterial genome assembly, as well as polishing tools for long-read assemblies, but many other options have been published. Most assemblers and polishing tools were run using default parameters, and it is possible that further optimization of settings for the individual software packages might have improved assemblies further than was demonstrated here. As sequencing technologies and assembly software continues to improve, continued validation of pipelines is advisable. Software such as poreTally provides user-friendly options for benchmarking genome assembly pipelines prior to implementation [62].

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

A secondary aim of this study was to identify plasmids in the hybrid assemblies. Automated tools have been developed

and validated for identification of plasmids from genome assemblies or read data, but they are dependant of collated databases of known plasmid sequences. As such, tools such as PlasmidFinder or mlplasmids can be applied for plasmid identification for *Enterobacteriaceae* or *Enterococcus faecium*, but *B. fragilis* is not supported at the time of writing [63, 64]. Therefore, we evaluated putative plasmid sequences by sequence identity and length comparison using the PLSDb webpage, identifying plasmid replication domains, and using circularization and relative coverage as indicators that a sequence represents a plasmid in a given isolate.

Only four of the twelve plasmid sequences from the seven isolates could be identified using the PLSDb and three of these were the same plasmid, pBFP35. Two other putative plasmids, pBFO18\_1 and pBFS01\_2, were likely plasmids pBF388c and pIP421 based on the partial sequences from these plasmids and plasmid length. This still leaves half of the circularized, putative plasmids unidentified. The two longer putative plasmids, pBFO17\_1 and pBFS01\_1, displayed a high degree of similarity, a mol% G+C out of the normal range for *B. fragilis* and a relative read depth of double the reads compared to the chromosome. Most annotated CDSs were associated with mobilizable elements, but no known plasmid replication domains could be identified. From the sequencing data alone, we cannot conclude that they represent true plasmids; however, the findings above and manual inspection of long-read mapping support that inference.

There are only 14 complete plasmid sequences from cultured *Bacteroides* isolates in the PLSDb v 2019\_03\_05, which is based on the NCBI RefSeq database. Many other *Bacteroides* plasmids have been partially described, and some are represented by partial sequences or marked as contig level in the NCBI nucleotide database [65–68]. Metagenomic sequencing and genome assembly projects are expanding the public sequence databases, and screening the NCBI nucleotide database, sequences with a high degree of similarity to the putative plasmid sequences from one patient isolate (BFO18) could be found. These originated from a rat caecum metagenomic plasmid sequencing project from Copenhagen, a few hours' drive from Odense University Hospital. To understand and perform surveillance of the dissemination of plasmids, there is a need for increased submissions of high quality, annotated and phenotypically validated sequences of bacterial isolates including plasmids. This study adds significantly to the number of complete plasmid sequences associated with *Bacteroides*.

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

We also intended to comprehensively identify resistance genes and IS elements in the hybrid genome assemblies. Using ABRicate with several resistance gene databases and IS-element nucleotide sequences, the findings of our previous study were confirmed and enhanced. Assemblies from Illumina sequencing alone would only allow partial IS

element identification [14]. With the complete assemblies, comprehensive identification of known IS elements upstream of the relevant resistance genes could be completed. In our first study, we used ResFinder with the available database at that time. Additionally, by including several databases, and lowering the % ID threshold, the number of genes identified increased. Lowering the % ID threshold resulted in identification of a possible *cfxA* allele in BFO17, putative *bexA* alleles in several isolates and nucleotide sequences with 66–71 % ID similarity to the *ugd* polymyxin-resistance gene but only 30–56 % COV. The later could possibly be the genes responsible for the inherent polymyxin resistance in *B. fragilis*. Genes with lower than 95 % nucleotide similarity can still encode proteins with 100 % amino acid similarity, as we found for the *bexB* alleles [14]. Therefore, lowering sequence ID thresholds can result in resistance gene alleles that are not (yet) included in nucleotide databases. Conversely, low thresholds can also result in false-positive identifications. Therefore, putative alleles should be analysed manually for proper interpretation. The defaults for running ABRicate (minimum % ID 75 and minimum % COV 0) are probably sufficient for most purposes.

As a result of the complete genome assembly of BFO17, we could now identify two copies of *nimJ*, while only one copy was identified in the short-read draft assembly of the same isolate in the previous study. Husain and colleagues identified the presence of three copies of *nimJ* in strain HMW615, when describing the *nimJ* gene [69]. We confirmed this finding by running ABRicate on the HMW615 assembly as done with the isolates of this study (not shown). Interestingly, RAST annotates a third *nim* gene (nucleotide positions 1 359 590...1 360 093) in the Unicycler hybrid assembly of BFO17, and the PGAP annotation includes an additional annotation of a pyridoxamine 5'-phosphate oxidase family gene (nucleotide positions 940 032...940 505), the family that includes the *nim* genes. It is possible that one or more novel alleles of the *nim* gene are present in BFO17.

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

*cfiA*, *nim* and *erm* genes were identified in isolates resistant to meropenem, metronidazole and clindamycin, respectively. However, there was not always co-occurrence of IS elements or identifiable known promoter sequences upstream of the resistance genes and phenotypical resistance. This suggests another unknown mechanism or promoter motifs that ensure sufficient expression of the genes to confer resistance, or the presence of other unknown elements that confer resistance in those isolates.

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

#### Funding information

Funding supporting this work was provided by a Danish Medical Research Grant [case no. 2013-5480/912523-108], as well as by internal funds from the Department of Clinical Microbiology, Odense University Hospital, Odense, Denmark. Thomas V. Sydenham's salary as a PhD student was funded by the University of Southern Denmark, and he received a travel grant from the Henrik and Emilie Ovesen Foundation.

#### Acknowledgements

We are very grateful to Professor Henrik Westh, Department of Clinical Microbiology, Hvidovre Hospital, Denmark, for allowing use of their Computerome High Performance Computing cluster account for data analysis, and to Mikala Wang, Department of Clinical Microbiology, Aarhus University Hospital, for gifting us the MDR *B. fragilis* strain isolated at her department. We thank Valentina Galata, Saarland University, Saarbrücken, Germany (first author of the paper describing PLSDB) and Natalie Ring, University of Bath, Bath, UK (first author of the 2018 reference by Ring and colleagues [18]) for kind and helpful answers to questions via e-mail.

#### Author contributions

The study was conceptualized by T. V. S. and U. S. J. Funding was secured by T. V. S., M. K., H. H. and U. S. J. Data curation and investigation was performed by T. V. S. Formal analysis was done by T. V. S. and S. O-P. 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. wrote the original draft and edited the manuscript. U. S. J., M. K., T. V. S., H. H., S. O-P. and H. W. revised the manuscript.

#### Conflicts of interest

The authors declare that there are no conflicts of interest

#### Ethical statement

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

#### Data bibliography

1. Sydenham TV *et al.*, Sequence read files [Oxford Nanopore Technologies (ONT) Fast5 files and Illumina fastq files], as well as the final genome assemblies, have been deposited in NCBI/ENA/DDJB under BioProject accession numbers PRJNA525024, PRJNA244942, PRJNA244943, PRJNA244944, PRJNA253771, PRJNA254401 and PRJNA254455 (2019).

2. Sydenham TV *et al.*, fastq format of demultiplexed ONT reads trimmed of adapters and barcode sequences, <https://doi.org/10.5281/zenodo.2677927> (2019).

3. Sydenham TV *et al.*, Genome assemblies from the assembly pipeline validation, <https://doi.org/10.5281/zenodo.2648546> (2019).



4. Cerdeño-Tárraga AM *et al.*, Reference assembly of *B. fragilis* CCUG4856<sup>T</sup>, [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_000025985.1/](https://www.ncbi.nlm.nih.gov/assembly/GCF_000025985.1/) (2005).

5. Sydenham TV *et al.*, Genome assemblies corresponding to each stage of the process of the assembly of the six MDR isolates, <https://doi.org/10.5281/zenodo.2661704>.

6. El-Gebali S *et al.*, Sequences of the plasmid replication domain families were downloaded from the Pfam database (<https://pfam.xfam.org/>). Specifically, the full-length sequences for all sequences in the full alignments were downloaded for: PF01051, PF01402, PF01446, PF01719, PF01815, PF02486, PF03090, PF03428, PF04796, PF05732, PF06504, PF06970, PF07042 and PF10134 (2019).

## References

- Wexler HM. Bacteroides: the good, the bad, and the nitty-gritty. *Clin Microbiol Rev* 2007;20:593–621.
- Nagy E, Urbán E, Nord CE, on behalf of ESCMID Study Group on Antimicrobial Resistance in Anaerobic Bacteria. Antimicrobial susceptibility of *Bacteroides fragilis* group isolates in Europe: 20 years of experience. *Clin Microbiol Infect* 2011;17:371–379.
- Ferlov-Schwensen SA, Sydenham TV, Hansen KCM, Hoegh SV, Justesen US. Prevalence of antimicrobial resistance and the *cfiA* resistance gene in Danish *Bacteroides fragilis* group isolates since 1973. *Int J Antimicrob Agents* 2017;50:552–.
- Nagy E, Justesen US, Eitel Z, Urbán E, ESCMID Study Group on Anaerobic Infection. Development of EUCAST disk diffusion method for susceptibility testing of the *Bacteroides fragilis* group isolates. *Anaerobe* 2015;31:65–71.
- Zankari E, Hasman H, Kaas RS, Seyfarth AM, Agersø Y *et al.* Genotyping using whole-genome sequencing is a realistic alternative to surveillance based on phenotypic antimicrobial susceptibility testing. *J Antimicrob Chemother* 2013;68:771–777.
- Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH *et al.* Predicting antimicrobial susceptibilities for *Escherichia coli* and *Klebsiella pneumoniae* isolates using whole genomic sequence data. *J Antimicrob Chemother* 2013;68:2234–2244.
- Boolchandani M, D'Souza AW, Dantas G. Sequencing-Based methods and resources to study antimicrobial resistance. *Nat Rev Genet* 2019;20:356–370.
- Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. Transforming clinical microbiology with bacterial genome sequencing. *Nat Rev Genet* 2012;13:601–612.
- Davies TJ, Stoesser N, Sheppard AE, Abuoun M, Fowler PW *et al.* Reconciling the potentially irreconcilable? Genotypic and phenotypic amoxicillin-clavulanate resistance in *Escherichia coli*. *BioRxiv* 2019;511402.
- Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M *et al.* The role of whole genome sequencing in antimicrobial susceptibility testing of bacteria: report from the EUCAST Subcommittee. *Clin Microbiol Infect* 2017;23:2–22.
- Nagy E, Becker S, Sóki J, Urbán E, Kostrzewa M. Differentiation of division I (*cfiA*-negative) and division II (*cfiA*-positive) *Bacteroides fragilis* strains by matrix-assisted laser desorption/ionization time-of-flight mass spectrometry. *J Med Microbiol* 2011;60:1584–1590.
- Rogers MB, Parker AC, Smith CJ. Cloning and characterization of the endogenous cephalosporinase gene, *cepA*, from *Bacteroides fragilis* reveals a new subgroup of Ambler class A beta-lactamases. *Antimicrob Agents Chemother* 1993;37:2391–2400.
- Rasmussen BA, Gluzman Y, Tally FP. Cloning and sequencing of the class B beta-lactamase gene (*ccrA*) from *Bacteroides fragilis* TAL3636. *Antimicrob Agents Chemother* 1990;34:1590–1592.
- Sydenham TV, Sóki J, Hasman H, Wang M, Justesen US *et al.* Identification of antimicrobial resistance genes in multidrug-resistant clinical *Bacteroides fragilis* isolates by whole genome shotgun sequencing. *Anaerobe* 2015;31:59–64.
- Ricker N, Qian H, Fulthorpe RR. The limitations of draft assemblies for understanding prokaryotic adaptation and evolution. *Genomics* 2012;100:167–175.
- Page AJ, De Silva N, Hunt M, Quail MA, Parkhill J *et al.* Robust high-throughput prokaryote de novo assembly and improvement pipeline for Illumina data. *Microb Genomics* 2016;2:mgen.0.000083.
- Schmid M, Frei D, Patrignani A, Schlapbach R, Frey JE *et al.* Pushing the limits of de novo genome assembly for complex prokaryotic genomes harboring very long, near identical repeats. *Nucleic Acids Res* 2018;46:8953–8965.
- Ring N, Abrahams JS, Jain M, Olsen H, Preston A *et al.* Resolving the complex *Bordetella pertussis* genome using barcoded nanopore sequencing. *Microb Genomics* 2018;4:mgen.0.000234.
- Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with multiplex MinION sequencing. *Microb Genom* 2017;3:mgen.0.000132.
- De Maio N, Shaw LP, Hubbard A, George S, Sanderson N *et al.* Comparison of long-read sequencing technologies in the hybrid assembly of complex bacterial genomes. *Microb Genom* 2019;5:mgen.0.000294.
- Cerdeño-Tárraga AM, Patrick S, Crossman LC, Blakely G, Abratt V *et al.* Extensive DNA inversions in the *B. fragilis* genome control variable gene expression. *Science* 2005;307:1463–1465.
- Kuwahara T, Yamashita A, Hirakawa H, Nakayama H, Toh H *et al.* Genomic analysis of *Bacteroides fragilis* reveals extensive DNA inversions regulating cell surface adaptation. *Proc Natl Acad Sci USA* 2004;101:14919–14924.
- Patrick S, Blakely GW, Houston S, Moore J, Abratt VR *et al.* Twenty-eight divergent polysaccharide loci specifying within- and amongst-strain capsule diversity in three strains of *Bacteroides fragilis*. *Microbiology* 2010;156:3255–3269.
- Nikitina AS, Kharlampieva DD, Babenko VV, Shirokov DA, Vakhitova MT *et al.* Complete genome sequence of an enterotoxigenic *Bacteroides fragilis* clinical isolate. *Genome Announc* 2015;3:e00450-15.
- Risse J, Thomson M, Patrick S, Blakely G, Koutsovoulos G *et al.* A single chromosome assembly of *Bacteroides fragilis* strain BE1 from Illumina and MinION nanopore sequencing data. *Gigascience* 2015;4:60.
- Soki J. *Bacteroides fragilis* S14 genome sequencing and assembly (data accessed on NCBI RefSeq database accession GCF\_001682215.1), [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_001682215.1/](https://www.ncbi.nlm.nih.gov/assembly/GCF_001682215.1/); 2015.
- Ho P-L, Yau C-Y, Wang Y, Chow K-H. Determination of the mutant-prevention concentration of imipenem for the two imipenem-susceptible *Bacteroides fragilis* strains, Q1F2 (*cfiA*-positive) and ATCC 25282 (*cfiA*-negative). *Int J Antimicrob Agents* 2018;51:270–271.
- Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 2018;34:i884–i890.
- Desai A, Marwah VS, Yadav A, Jha V, Dhaygude K *et al.* Identification of optimum sequencing depth especially for de novo genome assembly of small genomes using next generation sequencing data. *PLoS One* 2013;8:e60204.
- Wick RR, Judd LM, Holt KE. Deepbiner: demultiplexing barcoded Oxford Nanopore reads with deep convolutional neural networks. *PLoS Comput Biol* 2018;14:e1006583.
- De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack: visualizing and processing long-read sequencing data. *Bioinformatics* 2018;34:2666–2669.
- Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 2013;29:1072–1075.
- Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* 2015;25:1043–1055.



34. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P et al. BUSCO applications from quality assessments to gene prediction and phylogenomics. *Mol Biol Evol* 2018;35:543–548.
35. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014;30:2068–2069.
36. Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci USA* 2009;106:19126–19131.
37. Clark SC, Egan R, Frazier PI, Wang Z. ALE: a generic assembly likelihood evaluation framework for assessing the accuracy of genome and metagenome assemblies. *Bioinformatics* 2013;29:435–443.
38. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of *de novo* genome assemblies. *Bioinformatics* 2015;31:3350–3352.
39. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;34:3094–3100.
40. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 2009;25:1754–1760.
41. Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA et al. Circlator: automated circularization of genome assemblies using long sequencing reads. *Genome Biol* 2015;16:294.
42. Tatusova T, DiCuccio M, Badretdin A, Chetvernin V, Nawrocki EP et al. NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res* 2016;44:6614–6624.
43. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S et al. Identification of acquired antimicrobial resistance genes. *J Antimicrob Chemother* 2012;67:2640–2644.
44. Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P et al. Card 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. *Nucleic Acids Res* 2017;45:D566–D573.
45. Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. ISfinder: the reference centre for bacterial insertion sequences. *Nucleic Acids Res* 2006;34:D32–D36.
46. Galata V, Fehlmann T, Backes C, Keller A. PLSDB: a resource of complete bacterial plasmids. *Nucleic Acids Res* 2019;47:D195–202.
47. Jørgensen TS, Xu Z, Hansen MA, Sørensen SJ, Hansen LH. Hundreds of circular novel plasmids and DNA elements identified in a rat cecum metagenome. *PLoS One* 2014;9:e87924.
48. Shah HN. The genus *Bacteroides* and related taxa. In: Balows A, Trüper HG, Dworkin M and Harder W, Schleifer K-H (editors). *The Prokaryotes: a Handbook on the Biology of Bacteria: Ecophysiology, Isolation, Identification, Applications*. New York, NY: Springer New York; 1992. pp. 3593–3607.
49. Sóki J, Wareham DW, Rátkai C, Aduse-Opoku J, Urbán E et al. Prevalence, nucleotide sequence and expression studies of two proteins of a 5.6kb, class III, *Bacteroides* plasmid frequently found in clinical isolates from European countries. *Plasmid* 2010;63:86–97.
50. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res* 2016;44:D7–D19.
51. Trinh S, Haggoud A, Reysset G, Sebald M. Plasmids pIP419 and pIP421 from *Bacteroides*: 5-nitroimidazole resistance genes and their upstream insertion sequence elements. *Microbiology* 1995;141:927–935.
52. Haggoud A, Trinh S, Mourni M, Reysset G. Genetic analysis of the minimal replicon of plasmid pIP417 and comparison with the other encoding 5-nitroimidazole resistance plasmids from *Bacteroides* spp. *Plasmid* 1995;34:132–143.
53. Sóki J, Gal M, Brazier JS, Rotimi VO, Urbán E et al. Molecular investigation of genetic elements contributing to metronidazole resistance in *Bacteroides* strains. *J Antimicrob Chemother* 2006;57:212–220.
54. Hartmeyer GN, Sóki J, Nagy E, Justesen US. Multidrug-resistant *Bacteroides fragilis* group on the rise in Europe? *J Med Microbiol* 2012;61:1784–1788.
55. Darling AE, Mau B, Perna NT. progressiveMauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS One* 2010;5:e11147.
56. Nishida H. Comparative analyses of base compositions, DNA sizes, and dinucleotide frequency profiles in archaeal and bacterial chromosomes and plasmids. *Int J Evol Biol* 2012;2012:342482.
57. Brettin T, Davis JJ, Disz T, Edwards RA, Gerdes S et al. RASTtk: a modular and extensible implementation of the RAST algorithm for building custom annotation pipelines and annotating batches of genomes. *Sci Rep* 2015;5:8365.
58. Bayley DP, Rocha ER, Smith CJ. Analysis of *cepA* and other *Bacteroides fragilis* genes reveals a unique promoter structure. *FEMS Microbiol Lett* 2000;193:149–154.
59. Lin Y, Yuan J, Kolmogorov M, Shen MW, Chaisson M et al. Assembly of long error-prone reads using de Bruijn graphs. *Proc Natl Acad Sci USA* 2016;113:E8396–E8405.
60. Kamath GM, Shomorony I, Xia F, Courtade TA, Tse DN. Hinge: long-read assembly achieves optimal repeat resolution. *Genome Res* 2017;27:747–756.
61. Payne A, Holmes N, Rakyan V, Loose M. BulkVis: a graphical viewer for Oxford nanopore bulk FAST5 files. *Bioinformatics* 2019;35:2193–2198.
62. de Lannoy C, Risse J, de Ridder D. poreTally: run and publish *de novo* nanopore assembler benchmarks. *Bioinformatics* 2019;35:2663–2664.
63. Carattoli A, Zankari E, García-Fernández A, Voldby Larsen M, Lund O et al. In silico detection and typing of plasmids using PlasmidFinder and plasmid multilocus sequence typing. *Antimicrob Agents Chemother* 2014;58:3895–3903.
64. Arredondo-Alonso S, Rogers MRC, Braat JC, Verschuuren TD, Top J et al. mlplasmids: a user-friendly tool to predict plasmid- and chromosome-derived sequences for single species. *Microb Genom* 2018;4:mgen.0.000224.
65. Nguyen M, Vedantam G. Mobile genetic elements in the genus *Bacteroides*, and their mechanism(s) of dissemination. *Mob Genet Elements* 2011;1:187–196.
66. Shkoporov AN, Khokhlova EV, Kulagina EV, Smeianov VV, Kuchmiy AA et al. Analysis of a novel 8.9kb cryptic plasmid from *Bacteroides uniformis*, its long-term stability and spread within human microbiota. *Plasmid* 2013;69:146–159.
67. McNulty NP, Wu M, Erickson AR, Pan C, Erickson BK et al. Effects of diet on resource utilization by a model human gut microbiota containing *Bacteroides cellulosilyticus* WH2, a symbiont with an extensive glycobiome. *PLoS Biol* 2013;11:e1001637.
68. Pierce JV, Bernstein HD. Genomic diversity of enterotoxigenic strains of *Bacteroides fragilis*. *PLoS One* 2016;11:e0158171.
69. Husain F, Veeranagouda Y, Hsi J, Meggersee R, Abratt V et al. Two multidrug-resistant clinical isolates of *Bacteroides fragilis* carry a novel metronidazole resistance *nim* gene (*nimJ*). *Antimicrob Agents Chemother* 2013;57:3767–3774.
70. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *bioRxiv* 2019;530972.
71. Li H. Minimap and miniasm: fast mapping and *de novo* assembly for noisy long sequences. *Bioinformatics* 2016;32:2103–2110.
72. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. *Nat Biotechnol* 2019;37:540–546.
73. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH et al. Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome Res* 2017;27:722–736.
74. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 2012;19:455–477.
75. Antipov D, Korobeynikov A, McLean JS, Pevzner PA. hybridSPAdes: an algorithm for hybrid assembly of short and long reads. *Bioinformatics* 2016;32:1009–1015.
76. Souvorov A, Agarwala R, Lipman DJ. SKESA: strategic *k*-mer extension for scrupulous assemblies. *Genome Biol* 2018;19:153.
77. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. *PLoS Comput Biol* 2017;13:e1005595.

78. Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled *de novo* using only nanopore sequencing data. *Nat Methods* 2015;12:733–735.
79. Vaser R, Sovic I, Nagarajan N, Sikic M. Fast and accurate *de novo* genome assembly from long uncorrected reads. *Genome Res* 2017;27:737–746.
80. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* 2014;9:e112963.
81. Krumsiek J, Arnold R, Rattei T. Gepard: a rapid and sensitive tool for creating dotplots on genome scale. *Bioinformatics* 2007;23:1026–1028.
82. Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer. *Bioinformatics* 2011;27:1009–1010.
83. Arndt D, Grant JR, Marcu A, Sajed T, Pon A *et al.* PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res* 2016;44:W16–W21.

**Five reasons to publish your next article with a Microbiology Society journal**

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

**Find out more and submit your article at [microbiologyresearch.org](http://microbiologyresearch.org).**