

SOFTWARE

Open Access



MicroPIPE: validating an end-to-end workflow for high-quality complete bacterial genome construction

Valentine Murigneux^{1†}, Leah W. Roberts^{2,3,4*†}, Brian M. Forde², Minh-Duy Phan⁵,
Nguyen Thi Khanh Nhu⁵, Adam D. Irwin^{2,3}, Patrick N. A. Harris^{2,7}, David L. Paterson², Mark A. Schembri⁵,
David M. Whiley^{2,3} and Scott A. Beatson^{5,6*}

Abstract

Background: Oxford Nanopore Technology (ONT) long-read sequencing has become a popular platform for microbial researchers due to the accessibility and affordability of its devices. However, easy and automated construction of high-quality bacterial genomes using nanopore reads remains challenging. Here we aimed to create a reproducible end-to-end bacterial genome assembly pipeline using ONT in combination with Illumina sequencing.

Results: We evaluated the performance of several popular tools used during genome reconstruction, including base-calling, filtering, assembly, and polishing. We also assessed overall genome accuracy using ONT both natively and with Illumina. All steps were validated using the high-quality complete reference genome for the *Escherichia coli* sequence type (ST)131 strain EC958. Software chosen at each stage were incorporated into our final pipeline, MicroPIPE.

Further validation of MicroPIPE was carried out using 11 additional ST131 *E. coli* isolates, which demonstrated that complete circularised chromosomes and plasmids could be achieved without manual intervention. Twelve publicly available Gram-negative and Gram-positive bacterial genomes (with available raw ONT data and matched complete genomes) were also assembled using MicroPIPE. We found that revised basecalling and updated assembly of the majority of these genomes resulted in improved accuracy compared to the current publicly available complete genomes.

Conclusions: MicroPIPE is built in modules using Singularity container images and the bioinformatics workflow manager Nextflow, allowing changes and adjustments to be made in response to future tool development. Overall, MicroPIPE provides an easy-access, end-to-end solution for attaining high-quality bacterial genomes. MicroPIPE is available at <https://github.com/BeatsonLab-MicrobialGenomics/micropipe>.

Keywords: Nanopore, ONT, Pipeline, Sequence, Bacteria, Assembly, Polishing

* Correspondence: leah@ebi.ac.uk; scott.beatson@uq.edu.au

[†]Valentine Murigneux and Leah W. Roberts contributed equally to this work.

²University of Queensland Centre for Clinical Research, Brisbane, Queensland, Australia

⁵School of Chemistry and Molecular Biosciences, The University of Queensland, Brisbane, Queensland, Australia

Full list of author information is available at the end of the article



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Background

Bacterial genome construction using short-read sequencing has historically been difficult, largely due to the abundance of repeat sequences which collapse during de novo assembly, resulting in breaks in contiguous sequence [1]. However, long-read sequencing technologies, such as Oxford Nanopore Technology (ONT) and Pacific Biosciences (PacBio), are able to traverse these repeats enabling complete bacterial genomes [2]. Long reads also present the opportunity to correctly place single nucleotide variants (SNVs), particularly across complex regions of the genome that require more genomic context than short reads can provide. The accessibility and affordability of the ONT MinION sequencing device has resulted in its widespread use globally, allowing researchers the autonomy to perform their own experiments much more rapidly compared to using external sequencing facilities [3]. However, bacterial genome construction continues to be problematic, especially for non-specialised researchers.

Numerous tools designed to address aspects of complete bacterial genome construction have been developed by both ONT and community users, however few pipelines exist that offer end-to-end construction of bacterial genomes. Currently, these include Katuali [4], CCBGpipe [5], ASA³P [6] and Bactopia [7]. Katuali is an ONT-developed assembly pipeline implemented in Snakemake. It offers the user flexibility in software choice and is well-documented, but provides limited rationale or validation of the provided software. Additionally, it performs exclusively long-read assembly with no short-read polishing included. While ASA³P and Bactopia are able to generate assemblies using nanopore data, overall these pipelines were not designed solely for de novo assembly and are more focused on reproducible and comprehensive downstream analysis. CCBGpipe is distributed via Docker and implements a series of python scripts to run Canu with Racon and Nanopolish. However, similar to Katuali, this pipeline also performs Nanopore-only assembly (without Illumina) and was designed using Canu version 1.6, which is now several releases behind the current version (v2.1.1).

Substitution errors in nanopore reads have improved dramatically over recent years, from read accuracies of 60% [8] to the currently reported 95% for 1D reads using R9.4.1 flow cells [9]. While this is approaching that of Illumina (99.9%) [10] and PacBio (99%) [11], single nucleotide insertion/deletion (indel) errors remain problematic [12, 13]. Improvements in base-calling software (e.g. that account for methylation) and the introduction of the R10 pore have reduced these artefacts, but polishing nanopore assemblies with Illumina data has been generally required to achieve the highest quality possible [14].

With the rapid pace of ONT progression, development of new software and pipelines, or reappraisal of existing ones, has become an ongoing necessity. This has prompted the need for appropriate validation sets, to assess (or reassess) the accuracy of results. While simulated datasets provide an initial assessment of a tool's ability, data generated from biological sources provide additional confidence in its real-world application, as has been developed previously using metagenomic communities [15, 16]. *Escherichia coli* sequence type (ST)131 represents a globally disseminated lineage that has been intensively studied as a result of its recent emergence, antibiotic resistance and link to human disease [17–19]. Extensive knowledge of both *E. coli* (as a species) and the ST131 lineage makes it an ideal dataset to use for software and pipeline validation. Additionally, the *E. coli* ST131 strain EC958 represents an extensively curated and highly accurate reference genome, having been sequenced on multiple occasions using PacBio, Illumina and 454 pyrosequencing [20].

Here we present our complete pipeline, MicroPIPE, for automated construction of high-quality bacterial genomes using software chosen by systematic comparison of the most popular tools currently available in the community. Validation of each pipeline stage was completed using the high-quality *E. coli* ST131 reference genome, EC958. Subsequent validation of the complete pipeline was performed using 11 previously characterised ST131 *E. coli* strains, for which completely assembled genomes were already available. Finally, we tested MicroPIPE on 12 other publicly available bacterial isolates that had both a complete genome and associated raw nanopore sequencing data available. In all cases, we show that high-quality bacterial reference genomes can be achieved using MicroPIPE.

Implementation

Public data

The EC958 complete genome was downloaded from NCBI (GenBank: HG941718.1, HG941719.1, HG941720.1) [20]. Illumina reads for 12 ST131 genomes and draft assemblies for 95 ST131 were accessed from [17]. Twelve publicly available complete genomes were also selected to test MicroPIPE, under the following criteria: (i) the raw nanopore sequencing files (fast5) were available, (ii) a complete genome was made available for the same strain and (iii) Illumina sequencing data were available for the same strain. These 12 genomes represented 7 species from both gram-positive and gram-negative bacteria with chromosome sizes between 1.8 Mbp – 5.6 Mbps. A complete list of data used is provided in Supplementary dataset 1.

Culture and DNA extraction

Twelve ST131 *E. coli* isolates (including EC958) were grown from single colonies in Lysogeny Broth (LB) at 37 °C overnight with 250 rpm shaking. The overnight cultures (1.5 mL) were then pelleted for DNA extraction using the Wizard Genomic DNA Purification Kit (Promega) following manufacturer's protocol with modifications. Briefly, the cell pellet was lysed following the protocol for Gram negative bacteria. RNA was removed by 1 h incubation at 37 °C with RNase and the lysate was then mixed with Protein Precipitation Solution by vortexing for 5 s at max speed using Vortex-Genie 2 with horizontal tube adapter (Scientific Industries). The DNA was precipitated using isopropanol and washed with 70% ethanol. The DNA pellet was air-dried and then rehydrated in 100 µl EB buffer (QIAGEN) by incubation at 65 °C for 1 h. The DNA was quantified using a Qubit fluorometer (ThermoFisher Scientific) and the DNA fragment size was estimated using agarose gel electrophoresis (0.5% agarose in TAE, 90 V, 1h30m).

Nanopore sequencing

DNA from 12 ST131 *E. coli* were multiplexed onto a single FLO-MIN106 flow cell using the rapid barcode sequencing kit (SQK-RBK004) as per manufacturer's recommendation with the following adjustments: the barcoded DNA was pooled without a concentration step using AMPure XP beads prior to sequencing. Read metrics for each isolate are given in Supplementary Table 1.

Pipeline tools and settings

Specific parameters and commands used to perform the following analyses are provided in full in Supplementary dataset 1. MicroPIPE v0.8 uses Guppy v3.4.3, while MicroPIPE v0.9 uses Guppy v3.6.1.

Basecalling

Reads were basecalled using Guppy (v3.4.3) "fast" and "high-accuracy" modes. Fast mode was evaluated using both GPU and CPU servers, while the "high-accuracy" mode was evaluated using only GPU as the time to completion for this mode became unfeasible when run using CPUs. Upon the release of Guppy v3.6.1, reads were re-basecalled using only the "high-accuracy" mode. Guppy versions (3.4.3 and 3.6.1) were tested using the methylation aware config file "dna_r9.4.1_450bps_modbases_dam-dcm-cpg_hac.cfg".

Demultiplexing

Demultiplexing was evaluated using Guppy_barcode (v3.4.3) and qcat (v1.0.1) on the "passed" (>Q7) fastq reads after basecalling with Guppy. Demultiplexing using the raw fast5 reads was evaluated using Deepbiner

(v0.2.0) [21]. Demultiplexed fast5 reads were subsequently basecalled with Guppy (v3.4.3).

Quality control

Barcodes and adapters were trimmed using Porechop (v0.2.3_seqan2.1.1) (<https://github.com/rrwick/Porechop>). Overall read quality metrics and basecalling statistics were extracted using PycoQC (v2.2.3) [22]. Read length and quality metrics per sample were extracted using NanoPlot (v1.26.1) [23]. Average percentage read accuracy was determined by mapping the basecalled reads to the reference genome EC958 using Minimap2 (v2.17-r954-dirty) [24] and computing reads accuracy using Nanoplot. Filtering was evaluated using two tools: Filtlong (v0.2.0) (<https://github.com/rrwick/Filtlong>) and Japsa (v1.9-01a) (<https://github.com/mdcao/japsa/>).

Assembly

Six assemblers were evaluated for long-read assembly only: Canu (v1.9) [25], Flye (v2.5) [26], Raven (v1.1.5) (<https://www.nature.com/articles/s43588-021-00073-4>) (<https://github.com/lbcb-sci/raven>), Redbean (v2.5) [27], Shasta (v0.4.0: config file optimised for Nanopore: <https://github.com/chanzuckerberg/shasta/blob/master/conf/Nanopore-Dec2019.conf>) [28] and Unicycler (v0.4.7 long-read only) [29]. Three hybrid-assembly tools were also evaluated, including SPAdes (v3.13.1) [30], Unicycler (v0.4.7) and MaSuRCA (v3.3.5) [31]. Long-read correction was performed using Canu (v1.9).

Polishing and quality assessment

Polishing of the draft assemblies was evaluated using long reads (ONT), short reads (Illumina), and a combination of both long and short reads. Long read polishing was performed using Racon (v1.4.9) [32] and Medaka (v0.10.0) (<https://nanoporetech.github.io/medaka/>) (4 iterations of Racon based on Minimap2 v2.17-r941 overlaps followed by one iteration of Medaka), Nanopolish (v0.11.1) [33] (1 iteration based on Minimap2 v2.17-r941 alignment) and NextPolish (v1.1.0) [34] (2 iterations). Raw Illumina reads were trimmed using Trimmomatic (v0.36) [35] with the following settings: ILLUMINA-CLIP:TruSeq3-PE-2.fa:2:30:10 SLIDINGWINDOW:4:20 MINLEN:30. Short read polishing was performed using NextPolish (v1.1.0) and Pilon (v1.23) [36] (both 2 rounds of polishing based on BWA MEM v0.7.17-r1188 alignments).

Circularity was checked using NUCmer (v3.1) [37] to perform self-alignments. For Flye, Canu and Unicycler, circularisation was determined by the assemblers themselves. For Canu, circularisation was also confirmed using NUCmer self-alignments results and contigs were trimmed to remove overlapping ends. Circularisation for

Raven and Shasta was confirmed using generated GFA files. For MaSuRCA, circularisation was confirmed using Nucmer self-alignments results. For SPAdes, the plasmids were manually checked for circularity and the overlapping ends were trimmed. For Redbean, circularisation of the contigs was confirmed by alignment to the reference EC958 genome using QUAST.

Final assemblies were assessed for quality by comparison to the complete EC958 genome using the assess_assembly tool from Pomoxis (v0.3) (<https://github.com/nanoporetech/pomoxis>) as well as DNAdiff (v1.3) [37] and QUAST (v5.0.2) [38] to detect errors, misassemblies, and determine overall nucleotide identity.

Compute resources

All results were produced using cloud-based nodes with 16vCPUs and 32GB RAM. For the GPU node, the GPU is a NVIDIA Tesla P40 24GB while the CPUs are 2x Intel Xeon Silver 4214 2.2G (12C/24 T, 9.6GT/s, 16.5 M Cache, Turbo, HT [85 W] DDR4–2400).

ST131 phylogeny

ParSNP (v1.5.2) [39] was used to create an ST131 phylogeny using the 12 ST131 *E. coli* assembled in this study in addition to 95 ST131 *E. coli* short-read assemblies from Petty and Ben Zakour et al. [17]. Recombination was removed using PhiPack [40], as implemented in ParSNP. To evaluate the accuracy of each assembly and polishing step, we included our 12 completely polished assemblies (long and short read), 12 unpolished assemblies, 12 long-read polished assemblies and 12 short-read polished assemblies. The tree was visualised using Figtree (<http://tree.bio.ed.ac.uk/software/figtree/>) and iTOL [41].

MEME methylation motif analysis

The 20 bps sequence (−10 to +10) around the 401 shared SNPs were extracted using BEDTools getfasta (v2.28.0–33-g0f45761e) [42]. MEME (v5.2.0) [43, 44] was used to identify enriched motifs within the sequences using the default parameters of the classic mode and allowing zero or one occurrence per sequence. The motif CC(T/A)GG was significantly enriched in 393 sequences with an E-value of 6.2e-758.

Results

Validation of pipeline stages by comparison to EC958 complete genome

The main goal of this study was to create a robust and easily applicable pipeline for the construction of high-quality bacterial genomes with minimal manual manipulations. To achieve this, we first evaluated the performance of commonly used software at each stage of bacterial genome construction using the high-quality

EC958 genome (Accession: HG941718) as our standard for final genome accuracy. Figure 1 shows a diagram of the whole workflow, indicating the software chosen for comparison at each stage. Nanopore reads for EC958 were generated on a multiplexed run of 12 using the rapid barcoding kit on an R9.4.1 flow cell.

Basecalling

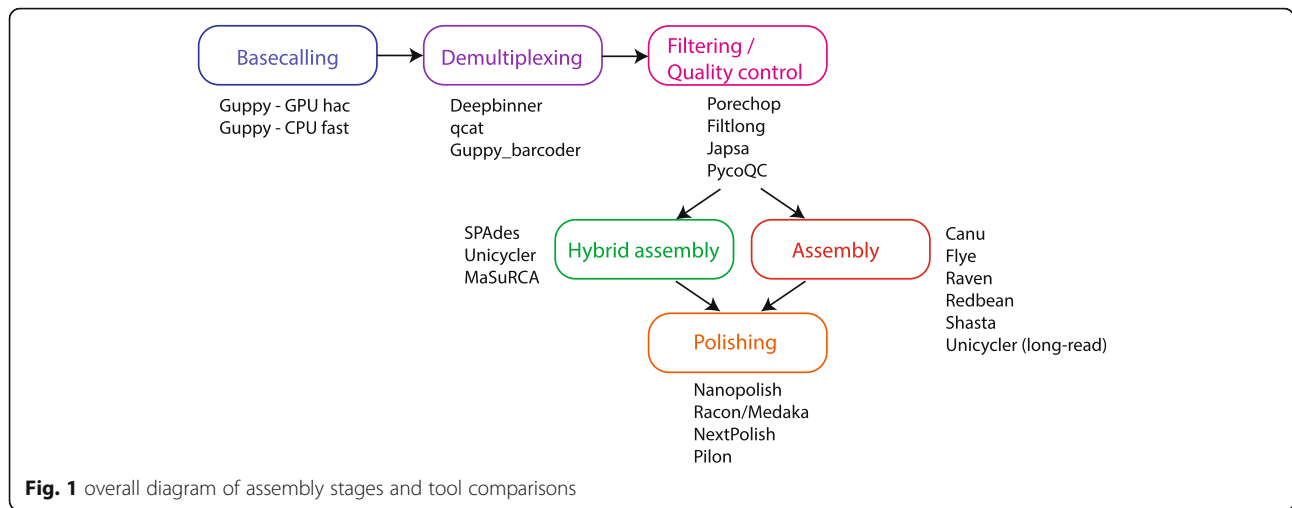
When considering the ongoing stability and accuracy of our overall pipeline, we decided to limit our basecalling validation to software that we were confident would be consistent and well-maintained for the foreseeable future. Many existing basecallers (such as Bonito, Flappie and Runnie) are currently research releases and therefore have minimal support and unknown longevity. Other basecallers are either depreciated (Albacore) or no longer updated (Scrappie). As such, we decided to focus our analysis on Guppy, which is the ONT recommended basecaller and is stably released and maintained.

Here we tested Guppy using both the “fast” and “high-accuracy” modes, as well as the CPU vs. GPU configurations. When using Guppy v3.4.3 with the “high-accuracy” setting on GPU servers we generated reads with approximately 91.0% accuracy in 828.5 min (13.81 h). Using the “fast” mode on CPUs, we were able to generate 88.9% accuracy in 2948.4 min (49.14 h) (Table 1). Testing the “high-accuracy” mode on a CPU server was unfeasible due to the time required for processing (fewer than 10% of reads completed basecalling in 1 week). Despite the lower per-read accuracy when using CPUs and the “fast” basecalling setting, the consensus quality of the overall finished genome (after assembly and polishing through MicroPIPE v0.8) was of comparable quality to that generated with the GPU and high-accuracy setting (Table 1).

We also tested the effects of methylation and found that using the “high-accuracy” model with methylation-aware basecalling achieved a similar per-read accuracy (90.6%) to the “high-accuracy” only model. The final assembly, however, had fewer SNPs (3 vs. 23 originally) and indels (31 vs. 45 originally) compared to the reference standard (Table 1).

Demultiplexing

For demultiplexing we tested three tools: Deepbiner [21], Guppy_barcode [45] and qcat [46]. While Guppy and qcat rely on basecalled reads, Deepbiner uses the raw fast5 reads. As such, we compared the total number of binned reads after both basecalling and binning for each tool. Overall, qcat was able to bin 89% of reads, compared to 84% for Guppy_barcode and 75% for Deepbiner (Supplementary Figure 1). Initially we chose qcat as the default demultiplexer as we prioritised read retention to maximise coverage of each genome.



However, following the recent depreciation of qcat (detailed on their GitHub: <https://github.com/nanoporetech/qcat>), ONT is recommending the use of the Guppy demultiplexer. As such, Guppy was chosen as the default demultiplexer for MicroPIPE, while qcat is still optionally available within the pipeline.

Filtering

Here we tested two filtering tools: Filtlong and Japsa. Filtlong has the advantage of being versatile enough to filter based on a number of requirements, such as read length, quality, percentage of reads to keep and the option of using an external reference. Japsa primarily filters

based on read length and quality. Read metrics after filtering using each tool are given in Supplementary Figure 2. Overall, we found that filtering with Japsa retained more reads, but with a reduced N50 read length and median read quality compared to Filtlong. Both tools took an equivalent amount of time to run. For all downstream analysis we filtered reads using Japsa with a minimum average quality cut-off of Q10 and 1 kb minimum read length, although Filtlong would have been equally suitable. Both filtering tools are available as optional steps in MicroPIPE. We have also included Rasusa [47] as an optional tool to randomly subsample large datasets down to a specific coverage (as necessary based

Table 1 Basecalling comparison: run-times, read accuracy and overall assembly accuracy

	Guppy3.4.3_hac	Guppy3.4.3_fast	Guppy3.4.3_hac_modbases	Guppy3.6.1_hac	Guppy3.6.1_hac_modbases
Basecalling comparison:					
Run time (ms)	49,707,952	176,906,144	57,479,661	57,977,178	46,296,565
Run time (h)	13.81	49.14	15.96	16.10	12.86
GPU/CPU	GPU	CPU	GPU	GPU	GPU
Num callers	4	16	8	8	8
Average read percent identity	91.0	88.9	90.6	93.7	91.0
Mean read quality	11.4	10.4	11.3	13.3	11.4
Number of binned reads (qcat)	240,766	233,802	238,847	244,830	240,156
Final assembly comparison:					
Assembly nucleotide identity (%)	99.99	99.99	99.99	99.99	99.99
Number of SNP (DNAdiff)	23	35	3	4	5
Number of indels (DNAdiff)	45	39	31	25	27
Assembly quality score (Pomoxis)	48.10	48.08	50.99	52.27	51.83
Mismatches per 100 kb (QUAST)	0.44	0.67	0.06	0.08	0.10
Indels per 100 kb (QUAST)	0.88	0.76	0.63	0.50	0.53

on user needs). This subsampling step is performed before trimming in order to reduce computational time.

Long-read-only assembly

A number of tools have been designed for de novo assembly from long reads. Here we compared six popular assembly tools and evaluated speed, completeness (of the chromosome and plasmids, including circularisation) and correctness (i.e. nucleotide identity) based on the complete EC958 reference genome standard, which contains 1 chromosome (5,109,767 bp) and 2 plasmids (135,602 bp and 4080 bp). Parameters used for all assemblers are given in Supplementary Dataset 1.

Overall, we found that all assemblers constructed the chromosome and larger (~ 135 kb) plasmid (Fig. 2, Supplementary Table 2). Raven, Redbean and Shasta did not assemble the smaller ~ 4 kb plasmid. While Canu was able to assemble both plasmids, closer inspection found them to be much larger than expected (1.4x and 2x larger for the large and small plasmid, respectively) due to overlapping ends that required additional trimming. Interestingly, both Flye and Canu assembled a third, previously unidentified, small plasmid of ~ 1.8 kb in size. This small plasmid was only identified when the Flye “--plasmids” mode was selected (to rescue short unassembled plasmids) and when certain or no filtering parameters were applied to the reads prior to assembly (Supplementary Table 3). Comparison of this small plasmid to the Illumina data for the EC958 reference genome standard confirmed its presence and was likely missed in the original assembly.

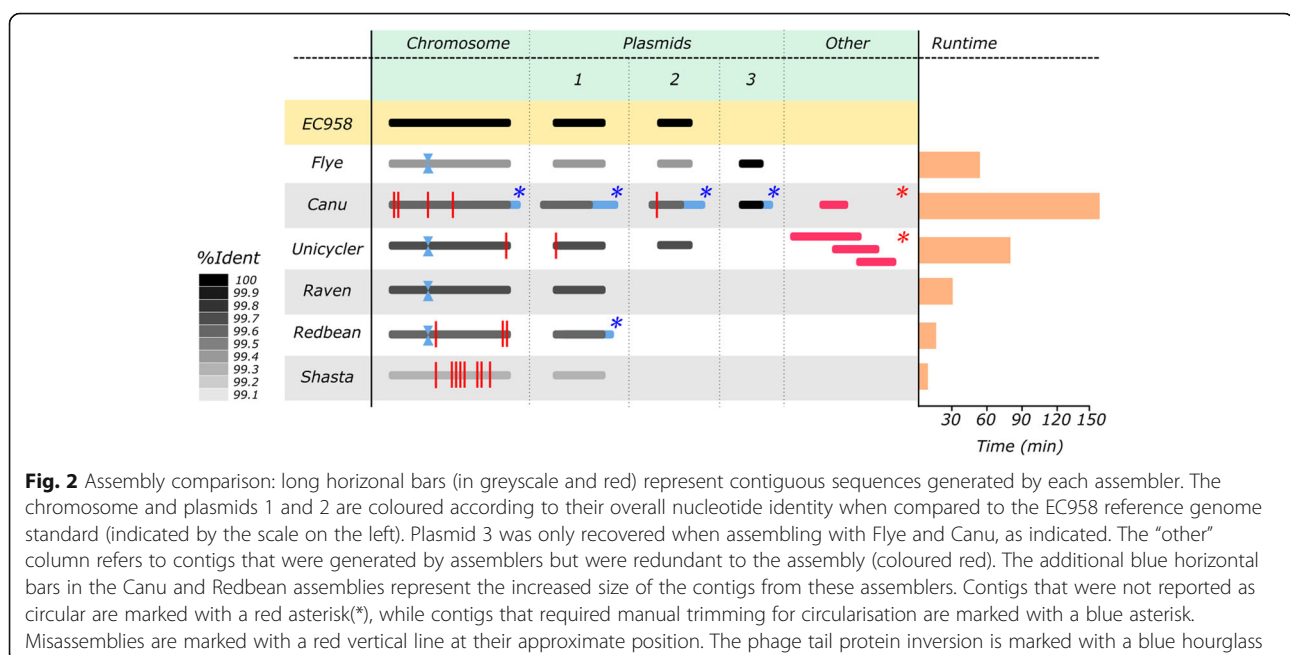
For most de novo assemblies, a number of small (< 4.5 kb) misassemblies were detected, mainly on the chromosome (Fig. 2). This included a small inversion, which on closer inspection was found to be an invertible phage tail protein that has been characterised previously [20]. This inversion was found in the Flye, Unicycler, Raven and Redbean assemblies and was not counted as a misassembly due to its biological relevance.

Additional contigs were found in both Canu and Unicycler (long-read only mode). The three additional contigs produced by Unicycler all matched other parts of the EC958 reference genome standard (two on the chromosome, one on the larger plasmid). The additional contig in Canu matched part of the additional ~ 1.8 kb plasmid.

In terms of speed, Shasta, Redbean and Raven were the fastest assemblers, completing in less than 30 min. Of the remainder, Flye was four times faster than Canu and two times faster than Unicycler. The majority of contigs from all assemblers were reported as circularised upon assembly completion, with the exception of the additional contigs in Canu and Unicycler. Redbean did not generate circularisation information, although the chromosome and plasmid contigs could be circularised manually or using 3rd party software following assembly. Overall, we found that Flye generated the best de novo assembly from long read data without the need for manual intervention.

Polishing

Polishing of assemblies generated using long reads is currently regarded as a necessity for ONT data due to



high per-read errors that can persist through to the de novo assemblies [14]. Here we tested the polishing capabilities of three different tools (Racon/Medaka, NextPolish and Nanopolish) using nanopore long reads against the de novo assembly generated using Flye. We additionally tested polishing with Illumina short reads (NextPolish and Pilon), which have a higher basecall accuracy. Polishing was tested both independently (i.e., long read and short read separately) as well as sequentially (long read followed by short read polishing) to determine the best polishing protocol.

Overall, we found that polishing with Racon and Medaka (four rounds of Racon and one round of Medaka, using long reads) followed by NextPolish (two rounds using short reads) achieved the most accurate assemblies (Fig. 3, Supplementary Table 4). Polishing using only long or short reads did not produce comparable levels of accuracy, therefore we emphasize the requirement of short read sequencing in parallel with Nanopore for high-quality complete genome assembly (as is already commonly done).

To confirm our choice of Flye as the best assembler, we polished assemblies generated from the other five long-read assemblers, described above, using this strategy (Supplementary Table 5). The polished Flye assembly remained the most accurate, closely followed by the polished Raven assembly.

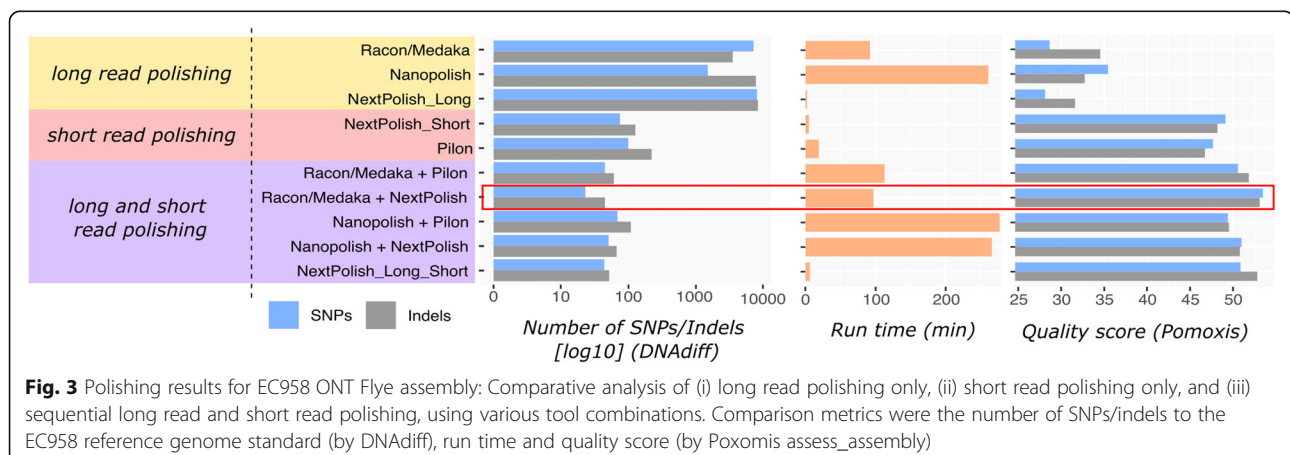
Hybrid assembly

In addition to long-read assembly (followed by short-read polishing), hybrid assemblers capable of using both long and short reads simultaneously have also been developed, and include Unicycler, MaSuRCA and SPAdes. Comparison of these pipelines to our genome completed with Flye, Racon, Medaka and NextPolish found that they did not outperform our current method. Unicycler was the only hybrid

assembler able to completely resolve the chromosome and both plasmids (SPAdes failed to circularise the chromosome while MaSuRCA was unable to assemble the 4 kb plasmid) (Supplementary Table 6). Additional long and short read polishing greatly improved the accuracy of the Unicycler and SPAdes hybrid assemblies but not MaSuRCA (Supplementary Table 5). We compared the quality of the genomes generated by either the best long-read only assembly (Flye) or the best hybrid assembler based on accuracy and structure (Unicycler) and polished with the same strategy. The polished assemblies contained a similar number of indels compared to the EC958 reference genome standard, however the Flye assembly contained around two-fold fewer substitution errors (Supplementary Table 5). Furthermore, Flye was nearly eight times faster than Unicycler (Supplementary Table 6).

Final pipeline

Based on the results of our comparative analysis for all of the major steps of bacterial genome assembly, we have developed MicroPIPE (Fig. 4). The pipeline is written in Nextflow [48] and the dependencies are packaged into Singularity [49] container images available through the Docker Hub and Quay.io BioContainers repositories. The bioinformatics workflow manager Nextflow allows users to run the pipeline locally or using common High-Performance Computing schedulers. Each step of the pipeline uses a specific container image which enables easy modifications to be made in the future to include new or updated tools. Furthermore, in addition to the recommended default pipeline settings, MicroPIPE also provides alternative software options and/or parameters to suit the user's individual needs. The pipeline is freely available on GitHub: <https://github.com/BeatsonLab/MicrobialGenomics/micropipe>.



Tool	Stage	Time				
		CPU	GPU	CPU (Total)	GPU (Total)	
Guppy	Basecalling	49.1 h (fast)	13.8 h (hac)			per MinION flow cell (12 <i>E. coli</i>)
Guppy	Demultiplexing	15 min	13.5 min	49.3 h	14.1 h	
pycoQC	Quality Control	1 min				
Porechop	Adapter trimming	30 min				per <i>E. coli</i>
Japsa	Filtering	5 min				
Flye	Assembly	38 min		120 min	98 min	
Racon + Medaka	Polishing (long reads)	42 min	20 min			
NextPolish	Polishing (short reads)	5 min				

Fig. 4 Overall pipeline: Stages and default tools in MicroPIPE. Stages in bold and italics are mandatory. All other pipeline steps are optional (users can start from fast5 or basecalled fastq files). Time for running each step is provided based on running 12 multiplexed *E. coli* samples with MicroPIPE v0.8. Basecalling (Guppy) and long-read polishing (Racon and Medaka) can be run on a GPU node. The rest of the pipeline is run using CPU resources. Fast = Guppy fast basecalling mode, hac = Guppy high accuracy basecalling mode. h = hour, min = minute

Evaluation of remaining differences with EC958 reference genome standard

The final genome for EC958 produced by MicroPIPE v0.8 was compared to the previously published EC958 reference genome standard (GenBank: HG941718.1) to assess any remaining differences. We observed a single 3.4 kb inversion corresponding to a phage tail protein switching event previously characterised in EC958 [20]. Overall, there were no other structural rearrangements. The final assembly contained an additional ~1.8 kb plasmid, with 100% nucleotide identity to previously reported *E. coli* plasmids (GenBank records CP048320.1, KJ484633.1, [50]). This plasmid appears to have been lost during size selection when constructing the original genomic DNA library for PacBio RSII sequencing of EC958 as it could be identified from de novo assembly of the corresponding Illumina reads.

Comparison of the two assemblies identified 68 remaining differences (66 on the chromosome, 2 on pEC958) (for full list, see Supplementary Dataset 1). The two differences in the plasmid sequence correspond to known errors in the EC958 reference genome standard (PacBio assembly constructed without Illumina polishing). The majority of the chromosomal differences were indels ($n = 45$, 67%) ranging from 1 to 6 bp in size. These

indels were mainly found in rRNA ($n = 31$), tRNA ($n = 4$), insertion sequences ($n = 4$), or phage-related genes ($n = 2$). The remaining 23 differences were SNPs, which were similarly found mainly in rRNA ($n = 13$) and insertion sequences ($n = 8$). These remaining differences likely represent an inability of current short-read polishing to adequately determine true alleles in repetitive regions of the genome. Using methylation-aware basecalling was found to significantly improve these errors, with only 3 SNPs and 31 indels (Supplementary Table 7).

MicroPIPE validation using 11 ST131 *E. coli*

To further test the robustness of MicroPIPE on other genomes, we included an additional 11 well-characterised *E. coli* ST131 strains [17] on a multiplexed run of 12 *E. coli* (i.e. 11 ST131 strains plus EC958).

Each strain took on average 120 min to run completely through MicroPIPE v0.8 using 16 threads (excluding the basecalling and demultiplexing steps) (Fig. 4). Of these 11 isolates, all had complete circularised chromosomes of the expected size. They also carried an array of plasmids, which were circularised in all cases except for a single isolate, HVM2044 (Supplementary Table 8). Re-analysis of this sample found that complete circularised plasmids can be achieved by adjusting the read filtering

step. We also identified additional small plasmids in six out of the 11 genomes ranging between 1.5–5 kb in size. Importantly, we found that these plasmids are not recovered when using filtering parameters above 1 kb.

In order to confirm the accuracy of the assemblies generated with MicroPIPE, we recreated the ST131 phylogeny from [17] using (i) the complete MicroPIPE assembly, (ii) long read only polished assembly, (iii) short read only polished assembly and (iv) unpolished Nanopore assembly, and assessed the position of each strain within the tree. We found that all MicroPIPE v0.8 assemblies and ONT assemblies polished with Illumina clustered closest to their Illumina counterpart within the phylogenetic tree (Fig. 5A). However, the long read polished and unpolished ONT assemblies in most cases did not cluster as expected. They also displayed longer branches indicative of the remaining errors within the assembly. Interestingly, the long read polished and unpolished assemblies for all ST131 isolates belonging to our previously defined fluoroquinolone-resistance clade

C [17, 18] clustered together independent of other clade C strains, possibly representing systematic errors from the ONT data. Further interrogation of the branch leading to this cluster identified 401 shared SNPs. Of these SNPs, 97% were transitions, particularly A → G ($n = 187$) and T → C ($n = 203$) (Supplementary Table 9, Fig. 5C). Further analysis of these sites determined that 393 (98%) were associated with a Dcm methylase motif CC(A/T)GG (Supplementary Figure 4).

MicroPIPE validation using publicly available ONT sequenced bacteria

Lastly, we tested MicroPIPE using 12 publicly available genomes from both gram-positive and gram-negative bacteria with available raw nanopore data (fast5) and validated our results using their corresponding complete genomes (Table 2, Supplementary Dataset 1). These genomes also represent a wide range of GC content to further validate the use of MicroPIPE on diverse bacterial species (Table 2). As most of these isolates were

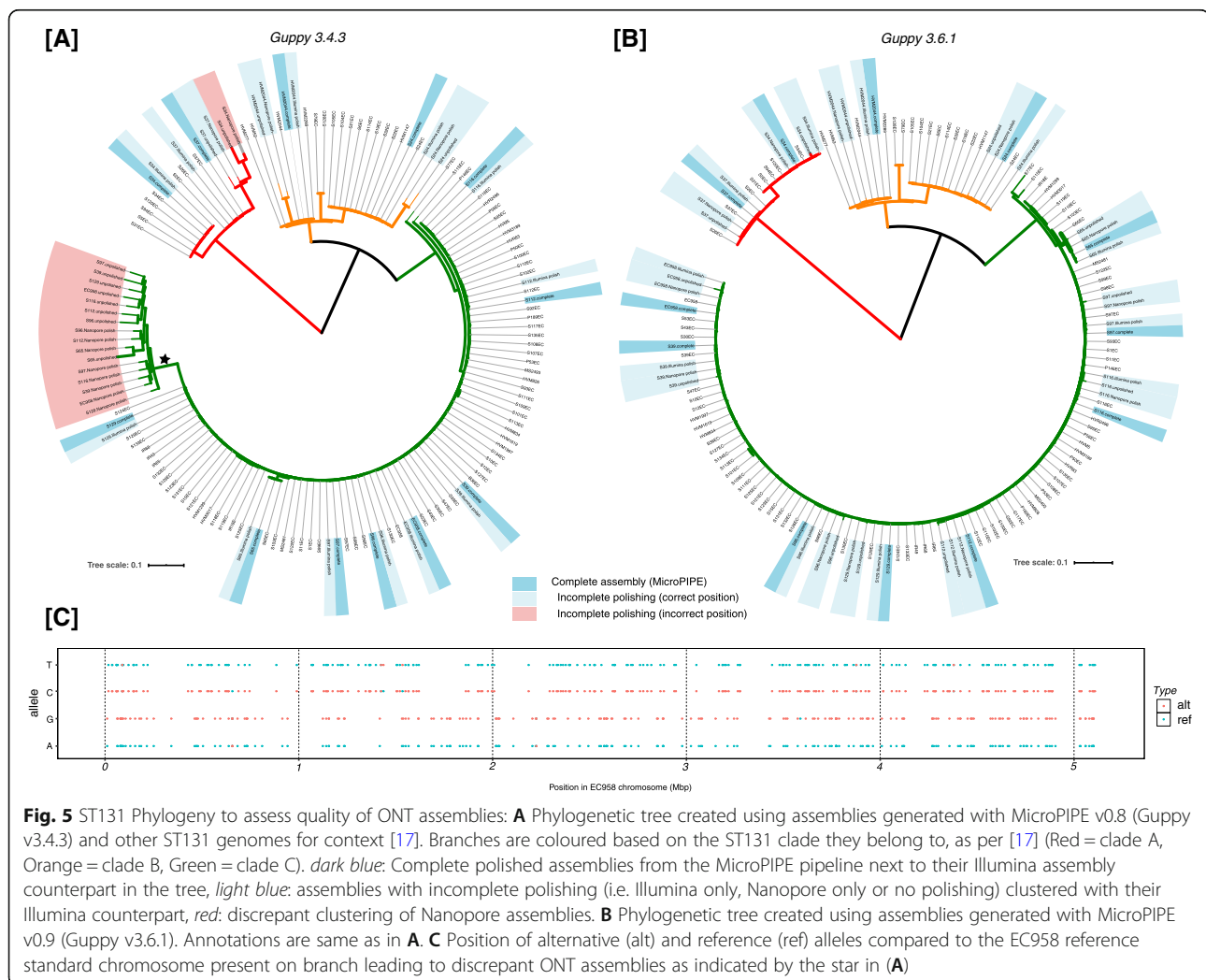


Table 2 MicroPIPE v0.9 results for public datasets

Reference	Strain	Reference genome assembly method and coverage	Chromosome/ plasmid	Reference genome size (bps)	Assembly size (bps)	GC content (%)	Circular?	Nucleotide identity (%)	DNAdiff SNPs GSNPs	DNAdiff Indels	QUAST misassemblies
Clement et al. [51]	<i>Salmonella enterica</i> serovar Napoli strain LC0541/17	Canu using Nanopore + Illumina 37x	Chromosome pLC0541_17	4,679,033 90,578	4,679,747 90,578	52.2	Yes	99.97	510 390	758	0
Sydenham et al. [52]	<i>Bacteroides fragilis</i> strain DCMOH0042B (BF042)	Unicycler using Nanopore + Illumina 200x	Chromosome pBF042_1 pBF042_2	5,141,257 83,06 5594	5,141,261 83,16 5629	43.3	Yes	99.99	65 9	14	0
Sydenham et al. [52]	<i>Bacteroides fragilis</i> strain CCLG48561	Unicycler using Nanopore + Illumina 200x	Chromosome pBF9343	5,205,133 36,560	5,205,138 36,559	43.1	Yes	99.99	25 5	22	1 (inversion)
Walker et al. [53]	<i>Streptococcus pyogenes</i> strain SP1336	Pacbio 105x	Chromosome	1,878,827	1,878,922	38.5	Yes	99.99	8 6	96	0
Wick et al. [14]	<i>Klebsiella pneumoniae</i> strain INF032	Unicycler using Nanopore + Illumina 133x	Chromosome	5,111,537	5,111,663	57.6	Yes	99.99	137 72	172	0
Taylor et al. [54]	<i>Escherichia coli</i> O157:H7 strain FS511705876	Unicycler using Nanopore + Illumina 692x	Chromosome pO157	5,483,434 94,581	5,483,452 94,593	50.4	Yes	99.99	52 2	103	0
Taylor et al. [54]	<i>Salmonella enterica</i> Bareilly strain CFSAN000189	Unicycler using Nanopore + Illumina 599x	Chromosome Plasmid	4,724,806 81,814	4,724,797 81,815	52.2	Yes	99.99	32 21	34	0
Bessonov et al. [55]	<i>Salmonella enterica</i> strain SA20055162	SMRT Analysis v. 1.3.3 using PacBio RS 80x	Chromosome Plasmid	4,730,612 78,193	4,640,715 105,679 98,127	51.7	Yes	99.99	0 53 13	15 22	3
Pitt et al. [56]	<i>Pandora fibrosis</i> strain 6399	Unicycler using Nanopore + Illumina 40x	Chromosome	5,592,065	5,592,075	62.8	Yes	99.99	28 5	9	0
Pitt et al. [56]	<i>Pandora fibrosis</i> strain 7641	Unicycler using Nanopore + Illumina 20x	Chromosome	5,592,064	5,591,941	62.8	Yes	99.99	81 15	102	0
Sieber et al. [57]	<i>Staphylococcus aureus</i> strain 110900	SPAdes using Nanopore + Illumina 334x	Chromosome Plasmid unnamed	2,918,239 2473	2,918,243 3356	32.7	Yes	99.99	6 6	5	0
Sieber et al. [57]	<i>Staphylococcus aureus</i> strain 128254	SPAdes using Nanopore + Illumina 219x	Chromosome Plasmid unnamed	2,877,083 2473	2,877,086	32.7	Yes	99.99	4 1	4	0

Flye was run using the --asm-coverage 100 parameter in order to reduce the computational run time. Only circular contigs are reported (as identified by Flye). For further details on all public data, see Supplementary dataset 1.

sequenced using entire flow cells, the coverage was reduced to 100x during the initial Flye assembly stage to minimise processing time.

Using MicroPIPE v0.9, we were able to completely assemble the chromosome and plasmids of all 12 isolates. We were also able to recover two additional plasmids from the *Salmonella enterica* str. SA20055162 that were not reported in the original assembly (Table 2).

To determine the accuracy of MicroPIPE, we compared our final assemblies with the submitted complete genome for each isolate. Overall, the fewest differences were detected between our MicroPIPE assembly and the complete genome of *Staphylococcus aureus* strain 110900 (6 SNPs, 5 indels) and strain 128254 (4 SNPs, 4 indels), constructed using ONT data basecalled with a recent version of Guppy (v3.2.6) (Table 2). These were followed by *Streptococcus pyogenes* strain SP1336, constructed using PacBio long-read sequencing (8 SNPs, 96 indels). All other comparisons yielded 25–510 SNPs, and 14–758 indels, with the greatest number of differences observed in the *Salmonella enterica* serovar Napoli strain LC0541/17 (Table 2).

With the exception of *S. pyogenes* SP1336, all other complete genomes were constructed using previously assembled nanopore data (Supplementary Dataset 1). Specifically, all assemblies with a high number of SNPs and indels were generated using reads basecalled with Albacore or a Guppy version prior to v3. As such, we hypothesise that our MicroPIPE assemblies likely represent corrections to the existing complete genomes, as a result of updated basecalling and assembly methods. Further investigation found that one sample, *Salmonella enterica* Bareilly str. CFSAN000189, also had a corresponding complete genome constructed using PacBio data. Comparison of our MicroPIPE assembly to this complete genome detected 0 SNPs and 15 indels, while there were 32 SNPs and 34 indels compared to the ONT complete genome.

Future development of MicroPIPE

Rapid and continual enhancement of nanopore technology has been integral to ONT's growth and popularity in recent years. It does, however, lead to several problems, including rapid depreciation, abandonment or replacement of software. As such, we have developed a modularised ONT/Illumina pipeline that can be readily adapted and re-evaluated alongside the changing nanopore landscape.

An example of MicroPIPE's adaptability came from the release of Guppy v3.6.1 during preparation of this manuscript. As this version reported a substantial increase in basecalling accuracy, we incorporated it into MicroPIPE (v0.9) and re-evaluated our pipeline's performance for all ST131 genomes.

Using MicroPIPE v0.9 on our EC958 data, we were able to resolve 21 out of the 23 SNPs and 32 out of 45 indels compared to the MicroPIPE v0.8 assembly (Guppy v3.4.3) (Supplementary Dataset 1, Supplementary Figure 3, Supplementary Table 7), relative to the published reference genome. Two SNPs and 12 indels were additionally detected using v3.6.1, which were not detected using v3.4.3. Both SNPs were detected in IS elements, while 11 out of the 12 indels were detected in rRNA genes. Overall, the v3.6.1 assembly performed better than the v3.4.3 assembly with only 29 differences compared to the complete reference EC958 genome (4 SNPs and 25 indels). Interestingly, using methylation-aware basecalling with Guppy v3.6.1 was not found to improve overall assembly accuracy (Supplementary Table 7).

We also found that by re-basecalling all other remaining ST131 isolates with MicroPIPE v0.9 and re-creating assemblies as before, we were able to achieve a remarkable increase in the accuracy of Nanopore-only assemblies, such that all assemblies clustered in their expected position within the tree (Fig. 5B).

Discussion

ONT long-read sequencing has quickly become one of the most prominent sequencing platforms for microbial researchers globally. However, despite the large number of bacterial genomes being completed using ONT, few end-to-end genome assembly pipelines exist. Here we created an easy, automated and reproducible genome assembly pipeline for the construction of complete, high-quality genomes using ONT in combination with Illumina sequencing. We also provide a robust, publicly available set of 12 ST131 genomes that can be used to validate future pipeline development or software advancements.

One of the main benefits of nanopore sequencing is its cost effectiveness, particularly when multiplexing several samples onto a single flow cell. Methods have been developed to improve yield and length during DNA extraction in order to achieve longer sequencing reads [15, 58]. However, here we show with our method that high-quality complete genomes can be achieved using a standard, commercially available DNA extraction kit coupled with up to 12 multiplexed samples. This builds on other advances such as those described by Wick et al. [59], and establishes an updated packaged pipeline that provides an efficient, cost effective and reproducible approach to bacterial genome construction.

In our comparative analysis of different aspects of bacterial genome assembly, we chose not to explore the effect of basecallers outside of ONT's Guppy_basecaller. As stated previously, many other existing basecallers have been released in a research-capacity (Bonito,

Flappie and Runnie), and are therefore unsuitable when considering the stability and maintenance of MicroPIPE. This is of particular importance to users from clinical settings, where consistency and versioning are essential when it comes to accrediting workflows. We were also confident that Guppy was among the highest performing basecallers, as this comparison has been completed previously [14]. Lastly, Guppy is ONTs recommended basecaller, coupled with several of Oxford Nanopore's devices, such as the MinIT, PromethION and GridION. For these reasons, we felt that it was in the best interest of the community at this time to provide a pipeline that used Guppy as the basecaller. We made a point of testing both the "high accuracy" mode on a GPU server compared to the "fast" mode on a CPU server, as not all Nanopore users are guaranteed to have access to GPU facilities. We found that, while the GPU server was significantly faster, basecalling reads using the "fast" mode with CPUs could also achieve high-quality genomes with MicroPIPE.

During preparation of this manuscript, Guppy v3.6.1 was released with a raw read accuracy of > 97% using R9.4.1 flow cells (<https://nanoporetech.com/accuracy>). Community feedback regarding this upgraded version supported increased overall accuracy, which prompted us to incorporate this version into our analysis (MicroPipe v0.9). We also found that Guppy v3.6.1 increased the overall accuracy of our assemblies, particularly where it came to unresolved indels using v3.4.3, which were suspected to be the result of technical artefacts around methylated sites [58]. Using Guppy v3.6.1 made Nanopore-only assemblies more feasible, particularly in cases where sufficient genetic context can be provided (e.g. identification of outbreak vs. non-outbreak strains). However, we found that overall both v3.4.3 and v3.6.1 still required polishing with short-read Illumina for maximum accuracy.

We observed some redundancy in the choice of tools for demultiplexing. Binning of reads with both Guppy_barcode and qcat performed almost equivalently (in terms of number of reads binned), with minimal differences in the overall assembly (Supplementary Table 10). Recent improvements to Guppy_barcode, which were released by ONT after compilation of this manuscript, suggest that Guppy_barcode is likely to be the default standard moving forward.

MicroPIPE implements a modest filtering measure to remove shorter, low quality reads from the dataset. In this study, we found that the length of sequencing reads used for assembly was an important parameter. Circularised chromosomes and large circularised plasmids were only obtained when the dataset contained a substantial proportion of reads longer than 5 kb (read length N50 for the 12 *E. coli* strains here ranged between 11 kb

and 15 kb). However, excessive removal of short reads negatively impacted the recovery of small plasmids, where removing reads ≤ 2 kb resulted in the loss of several small plasmids in a number of strains (data not shown). This was also the case when using certain additional filtering parameters with Filtlong, where "--min-length 1000 --keep_percent 90" resulted in the loss of the ~ 1.8 kb small plasmid identified in EC958, which was retained when filtering with Japsa at "--min-length 1000" (Supplementary Tables 11 and 12). As such, we have implemented a conservative 1 kb filtering cut-off (using Japsa) as default in MicroPIPE to retain reads and small plasmids.

We also found when testing MicroPIPE on publicly available data that harsher filtering is sometimes desirable, especially in cases where a single bacterial genome has been sequenced using an entire flow cell (such that we used the Flye parameter "--asm-coverage 100" to reduce coverage for initial disjoint assembly). As such, pre-processing of large quantities of highly ununiform data may be the most desirable method. This is possible to implement within MicroPIPE, as users may choose to randomly subsample with Rasusa, or implement subsampling and filtering with Filtlong over Japsa (the current default tool). Ultimately, understanding the quality and read lengths of the input data is a valuable step in generating the best possible assembly. We also provided the user read quality assessment using PycoQC to assist in parameter selection.

Several other comparative analyses have been published exploring the overall utility of different assemblers, in particular Wick et al. [60], who provide a comprehensive assembly comparison using both simulated and real read datasets. While we did not test NECAT and Miniasm, we found that our results generally matched those reported by Wick et al., particularly when it came to the overall strong performance of Flye. The most recent version of Flye (v2.8) also removes the need to nominate a genome size, making it a more robust option. However, we found that this version did not outperform the release used in this paper (v2.5) on our dataset, as it was unable to circularise all plasmids. As such, we have retained Flye v2.5 in MicroPIPE.

Long and short read polishing is a staple of high-quality genome assembly, as the combination of both ensures the correct contextual placement of variants as well as highly accurate basecalls. However, while long-reads have enabled completion of assemblies by spanning repetitive regions, polishing of these regions with short reads remains a problem. Here we found that the majority of remaining differences between our EC958 ONT assembly and the reference assembly (constructed with PacBio single molecule real time [SMRT] sequencing) resided in repetitive regions. Ideally, polishing with

long reads only would be a viable method to reduce these errors as they would have sufficient coverage to ensure correct placement of the repeat variant. However, as we show here, long read-only polishing was insufficient (likely due to per-read accuracy), and short read polishing was necessary for removal of the majority of errors. Currently, final polishing and assembly prior to completion will still necessitate manual frameshift inspection. While impractical and costly, a combination of both PacBio and ONT assembly could correct inherent biases in both technologies, using a consensus tool such as Trycycler (<https://github.com/rrwick/Trycycler>).

Long-read correction could also provide another means of error reduction [61, 62]. Upon subsequent analysis, we did find that the final assembly produced when using MicroPIPE v0.9 with Canu error-corrected reads was marginally better than using raw reads (Supplementary Table 13). However, this was at a cost of a 2.5 times slower runtime. We further tested raw vs. corrected reads with the latest Guppy version (v4.4.1) and did not see improvement of the final genome with corrected reads. Additionally, Flye (as the default assembler in MicroPIPE) recommends the use of raw reads over corrected reads (<https://github.com/fenderglass/Flye/blob/flye/docs/USAGE.md#error-corrected-reads-input>). As such, we have not implemented read-correction in our pipeline, but it could be implemented by users separately (or added in their own version of the pipeline) if desired.

We validated MicroPIPE using a set of 12 well-characterised *E. coli* isolates described previously from a global collection [17, 18]. We did this for several reasons, including (i) the availability of an existing high-quality reference genome and associated phylogenetic data (ii) the robustness of *E. coli* as a representative species and workhorse organism, and (iii) our extensive knowledge of the *E. coli* genome and ST131 lineage. We hope that by providing this dataset to the wider community, it can serve as a resource for future validation and testing of not only MicroPIPE, but other microbial assembly pipelines and tools.

In addition to in-house ONT sequencing data, we also tested MicroPIPE on a variety of publicly available bacterial genomes to evaluate its assembly capabilities on other species. Without any manual intervention, MicroPIPE was able to assemble all 12 genomes, while also recovering additional plasmids that were likely missed in the original assembly. When evaluating correctness of the genomes, we found a number of remaining SNPs and indels when compared to the complete genomes provided. Investigation into construction of the reference genomes found that 11 of the 12 genomes provided were constructed previously using ONT sequencing data, leading us to believe that differences in our assemblies compared to the “reference” genomes may actually be

corrections. Indeed, the genomes with the closest match between reference and MicroPIPE assembly were the genomes constructed using PacBio or ONT data with contemporary basecalling. As such, we believe that genomes completed historically using ONT reads should be used cautiously, and raw ONT data provided where possible to allow for reconstruction and improvement of the assembly as the technology improves.

Conclusions

Overall, we present an end-to-end pipeline for high-quality bacterial genome construction designed to be easily implemented in the research lab setting. We believe this will be a useful resource for users to easily and reproducibly construct complete bacterial genomes from Nanopore sequencing data.

Availability and requirements

Project name: MicroPIPE

Project home page: <https://github.com/BeatsonLab-MicrobialGenomics/micropipe>

Operating system(s): Linux/Unix/Mac

Programming language: Nextflow, Python

Other requirements: Java 8 or higher, Singularity 2.3.x or higher, Oxford Nanopore Technologies community access (Guppy)

License: GNU GPL-v3

Any restrictions to use by non-academics: None

Abbreviations

ONT: Oxford Nanopore Technology; PacBio: Pacific Biosciences; ST: Sequence type; bp: Base pair; kb: Kilobase pair; SNP: Single nucleotide polymorphism; indel: Insertion/deletion; IS: Insertion sequence

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12864-021-07767-z>.

Additional file 1: Table S1. Read length and quality metrics per isolate. **Table S2.** assembly results for EC958. **Table S3.** EC958 assembly results using different Flye parameters, demultiplexing tools and read filtering parameters. **Table S4.** Polishing tool comparison **Table S5.** Polishing comparison using other assemblers **Table S6.** Hybrid assembly comparison to Flye+Racon/medaka+NextPolish **Table S7.** final assembly comparisons between Guppy versions and methylation-aware basecalling. **Table S8.** MicroPIPE v0.8 results for 11 ST131. **Table S9.** SNP types for clade C unpolished/Illumina unpolished assemblies. **Table S10.** Demultiplexing comparison between qcat and Guppy **Table S11.** Assembly comparison using different filtering parameters(qcat demultiplexing). **Table S12.** Assembly comparison using different filtering parameters (guppy demultiplexing). **Table S13.** Assembly comparison using ONT raw or corrected reads **Figure S1.** Demultiplexing metrics. **Figure S2.** Long-read metrics using different demultiplexing tools and read filtering parameters (using EC958 ONT data). **Figure S3.** comparison of SNPs/indels in ONT assemblies vs. complete EC958 chromosome. **Figure S4.** Motif enriched in the sequences around the 401 shared SNPs from the branch leading to discrepant ONT assemblies.

Additional file 2. Strain list 1: ST131 data used for validation (including accessions for complete genomes, raw ONT data and Illumina read data). **Strain list 2:** Accession and metadata for the eight publicly

available ONT assemblies used in subsequent validation. **Execution parameters:** Specific commands used for each tool in this study. **Guppy 3.4.3 vs. Guppy 3.6.1 remaining variants:** List of variant positions and genetic context for Guppy v3.4.3 and v3.6.1 based on comparison to the EC958 reference genome standard. **MicroPIPE metrics:** Pipeline run times and memory metrics for 12 ST131 strains.

Acknowledgements

We would like to acknowledge Thom Cuddihy (QCIF Facility for Advanced Bioinformatics) for his assistance and advice regarding pipeline testing and development. This research was supported by QRIScloud and by use of the Nectar Research Cloud. The Nectar Research Cloud is a collaborative Australian research platform supported by the National Collaborative Research Infrastructure Strategy (NCRIS).

Authors' contributions

All authors (VM, LWR, BF, MDP, NTKN, ADI, PNAH, DLP, MAS, DMW, SAB) conceptualised the study. VM, LWR, BMF and SAB developed the methodology. MDP and MAS provided the bacterial strains and ONT sequencing data. VM wrote the pipeline. VM, LWR and NTKN conducted formal analysis. All authors (VM, LWR, BF, MDP, NTKN, ADI, PNAH, DLP, MAS, DMW, SAB) contributed to the interpretation of results. SAB and MAS supervised aspects of the project and provided essential expert analysis. LWR and VM wrote the original manuscript. BMF, SAB and MAS edited the manuscript. All authors read and approved the final manuscript.

Funding

LWR was supported by a Sakzewski Translational Research Grant. This work was supported by funding from the Queensland Genomics Health Alliance (now Queensland Genomics), Queensland Health, the Queensland Government.

Availability of data and materials

The datasets generated and analysed during the current study are available under the following Bioprojects (specific accessions available in supplementary dataset 1): EC958 complete genome (GenBank: HG941718.1), ST131 Illumina data (PRJEB2968), ST131 Nanopore data (fast5 and fastq [demultiplexed]; PRJNA679678).

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

None to declare.

Author details

¹QCIF Facility for Advanced Bioinformatics, Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland, Australia. ²University of Queensland Centre for Clinical Research, Brisbane, Queensland, Australia. ³Queensland Children's Hospital, Brisbane, Queensland, Australia. ⁴European Bioinformatics Institute, European Molecular Biology Laboratory (EMBL), Hinxton, Cambridge, UK. ⁵School of Chemistry and Molecular Biosciences, The University of Queensland, Brisbane, Queensland, Australia. ⁶Australian Centre for Ecogenomics, The University of Queensland, Brisbane, Queensland, Australia. ⁷Central Microbiology, Pathology Queensland, Royal Brisbane & Women's Hospital, Brisbane, Queensland, Australia.

Received: 9 March 2021 Accepted: 3 June 2021

Published online: 25 June 2021

References

- Klassen JL, Currie CR. Gene fragmentation in bacterial draft genomes: extent, consequences and mitigation. *BMC Genomics*. 2012;13(1):14. <https://doi.org/10.1186/1471-2164-13-14>.

- Koren S, Phillippy AM. One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly. *Curr Opin Microbiol*. 2015;23:110–20. <https://doi.org/10.1016/j.mib.2014.11.014>.
- Lemon JK, Khil PP, Frank KM, Dekker JP. Rapid nanopore sequencing of plasmids and resistance gene detection in clinical isolates. *J Clin Microbiol*. 2017;55(12):3530–43. <https://doi.org/10.1128/JCM.01069-17>.
- Katuali. ONT (Oxford Nanopore Technology); 2020. <https://github.com/nanoporetech/katuali>. Accessed Apr 2021.
- Liao YC, Cheng HW, Wu HC, Kuo SC, Lauderdale TY, Chen FJ. Completing circular bacterial genomes with assembly complexity by using a sampling strategy from a single MinION run with barcoding. *Front Microbiol*. 2019;10:2068. <https://doi.org/10.3389/fmicb.2019.02068>.
- Schwengers O, Hoek A, Fritzenwanker M, Falgenhauer L, Hain T, Chakraborty T, et al. ASA3P: an automatic and scalable pipeline for the assembly, annotation and higher-level analysis of closely related bacterial isolates. *PLoS Comput Biol*. 2020;16(3):e1007134. <https://doi.org/10.1371/journal.pcbi.1007134>.
- Petit RA 3rd, Read TD. Bactopia: a flexible pipeline for complete analysis of bacterial genomes. *mSystems*. 2020;5(4):e00190.
- Rang FJ, Kloosterman WP, de Ridder J. From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy. *Genome Biol*. 2018;19(1):90. <https://doi.org/10.1186/s13059-018-1462-9>.
- R10.3: the newest nanopore for high accuracy nanopore sequencing – now available in store. <https://nanoporetech.com/about-us/news/r103-newest-nanopore-high-accuracy-nanopore-sequencing-now-available-store>. Accessed Apr 2021.
- Measuring sequencing accuracy. <https://emea.illumina.com/science/technology/next-generation-sequencing/plan-experiments/quality-scores.html>. Accessed Apr 2021.
- Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, et al. Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat Biotechnol*. 2019;37(10):1155–62. <https://doi.org/10.1038/s41587-019-0217-9>.
- Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. Opportunities and challenges in long-read sequencing data analysis. *Genome Biol*. 2020;21(1):30. <https://doi.org/10.1186/s13059-020-1935-5>.
- Weirather JL, de Cesare M, Wang Y, Piazza P, Sebastiano V, Wang XJ, et al. Comprehensive comparison of Pacific Biosciences and Oxford Nanopore Technologies and their applications to transcriptome analysis. *F1000Res*. 2017;6:100.
- Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for Oxford Nanopore sequencing. *Genome Biol*. 2019;20(1):129. <https://doi.org/10.1186/s13059-019-1727-y>.
- Nicholls SM, Quick JC, Tang S, Loman NJ. Ultra-deep, long-read nanopore sequencing of mock microbial community standards. *Gigascience*. 2019;8(5):giz043.
- Sevim V, Lee J, Egan R, Clum A, Hundley H, Lee J, et al. Shotgun metagenome data of a defined mock community using Oxford Nanopore, PacBio and Illumina technologies. *Sci Data*. 2019;6(1):285. <https://doi.org/10.1038/s41597-019-0287-z>.
- Petty NK, Ben Zakour NL, Stanton-Cook M, Skippington E, Totsika M, Forde BM, et al. Global dissemination of a multidrug resistant *Escherichia coli* clone. *Proc Natl Acad Sci U S A*. 2014;111(15):5694–9. <https://doi.org/10.1073/pnas.1322678111>.
- Ben Zakour NL, Alsheikh-Hussain AS, Ashcroft MM, Khanh Nhu NT, Roberts LW, Stanton-Cook M, et al. Sequential acquisition of virulence and fluoroquinolone resistance has shaped the evolution of *Escherichia coli* ST131. *mBio*. 2016;7(2):e00347–16.
- Johnson JR, Porter S, Thuras P, Castanheira M. The pandemic H30 subclone of sequence type 131 (ST131) as the leading cause of multidrug-resistant *Escherichia coli* infections in the United States (2011–2012). *Open Forum Infect Dis*. 2017;4(2):ofx089.
- Forde BM, Ben Zakour NL, Stanton-Cook M, Phan MD, Totsika M, Peters KM, et al. The complete genome sequence of *Escherichia coli* EC958: a high quality reference sequence for the globally disseminated multidrug resistant *E. coli* O25b:H4-ST131 clone. *PLoS One*. 2014;9(8):e104400.
- Wick RR, Judd LM, Holt KE. Deepbinner: Demultiplexing barcoded Oxford Nanopore reads with deep convolutional neural networks. *PLoS Comput Biol*. 2018;14(11):e1006583. <https://doi.org/10.1371/journal.pcbi.1006583>.
- Leger A, Leonardi T. pycoQC, interactive quality control for Oxford Nanopore Sequencing. *J Open Source Softw*. 2019;4(34):1236.

23. De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack: visualizing and processing long-read sequencing data. *Bioinformatics*. 2018; 34(15):2666–9. <https://doi.org/10.1093/bioinformatics/bty149>.
24. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34(18):3094–100. <https://doi.org/10.1093/bioinformatics/bty191>.
25. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*. 2017;27(5):722–36. <https://doi.org/10.1101/gr.215087.116>.
26. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. *Nat Biotechnol*. 2019;37(5):540–6. <https://doi.org/10.1038/s41587-019-0072-8>.
27. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods*. 2020;17(2):155–8. <https://doi.org/10.1038/s41592-019-0669-3>.
28. Shafin K, Pesout T, Lorig-Roach R, Haukness M, Olsen HE, Bosworth C, et al. Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes. *Nat Biotechnol*. 2020;38(9):1044–53. <https://doi.org/10.1038/s41587-020-0503-6>.
29. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. *PLoS Comput Biol*. 2017; 13(6):e1005595. <https://doi.org/10.1371/journal.pcbi.1005595>.
30. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol*. 2012;19(5):455–77. <https://doi.org/10.1089/cmb.2012.0021>.
31. Zimin AV, Puiu D, Luo MC, Zhu T, Koren S, Marçais G, et al. Hybrid assembly of the large and highly repetitive genome of *Aegilops tauschii*, a progenitor of bread wheat, with the MaSuRCA mega-reads algorithm. *Genome Res*. 2017;27(5):787–92. <https://doi.org/10.1101/gr.213405.116>.
32. Vaser R, Sovic I, Nagarajan N, Sikic M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res*. 2017;27(5):737–46. <https://doi.org/10.1101/gr.214270.116>.
33. Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat Methods*. 2015;12(8): 733–5. <https://doi.org/10.1038/nmeth.3444>.
34. Hu J, Fan J, Sun Z, Liu S. NextPolish: a fast and efficient genome polishing tool for long-read assembly. *Bioinformatics*. 2020;36(7):2253–5. <https://doi.org/10.1093/bioinformatics/btz891>.
35. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014;30(15):2114–20. <https://doi.org/10.1093/bioinformatics/btu170>.
36. Walker BJ, Abeel T, Shea T, Priest M, Bouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One*. 2014;9(11):e112963. <https://doi.org/10.1371/journal.pone.0112963>.
37. Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. MUMmer4: a fast and versatile genome alignment system. *PLoS Comput Biol*. 2018;14(1):e1005944. <https://doi.org/10.1371/journal.pcbi.1005944>.
38. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics*. 2013;29(8):1072–5. <https://doi.org/10.1093/bioinformatics/btt086>.
39. Treangen TJ, Ondov BD, Koren S, Phillippy AM. The harvest suite for rapid core-genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome Biol*. 2014;15(11):524. <https://doi.org/10.1186/s13059-014-0524-x>.
40. Bruen TC, Philippe H, Bryant D. A simple and robust statistical test for detecting the presence of recombination. *Genetics*. 2006;172(4):2665–81. <https://doi.org/10.1534/genetics.105.048975>.
41. Letunic I, Bork P. Interactive tree of life (iTOL) v4: recent updates and new developments. *Nucleic Acids Res*. 2019;47(W1):W256–9. <https://doi.org/10.1093/nar/gkz239>.
42. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010;26(6):841–2. <https://doi.org/10.1093/bioinformatics/btq033>.
43. Bailey TL. Discovering novel sequence motifs with MEME. *Curr Protoc Bioinformatics*. 2002;Chapter 2:Unit 2.4.
44. Bailey TL, Elkan C. Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *Proc Int Conf Intell Syst Mol Biol*. 1994;2:28–36.
45. Guppy Barcoder. Oxford Nanopore Technology; 2020. https://community.nanoporetech.com/protocols/Guppy-protocol/v/gpb_2003_v1_rev_w_14dec2018/barcoding-demultiplexing. Accessed Apr 2021.
46. qcat demultiplexer. ONT (Oxford Nanopore Technology); 2020. <https://github.com/nanoporetech/qcat>. Accessed 17 May 2019.
47. Hall M. mbhall88/rasusa 0.3.0 (Version 0.3.0). Zenodo; 2020. <https://doi.org/10.5281/zenodo.3731394>. Accessed Jan 2021.
48. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow enables reproducible computational workflows. *Nat Biotechnol*. 2017;35(4):316–9. <https://doi.org/10.1038/nbt.3820>.
49. Kurtzer GM, Sochat V, Bauer MW. Singularity: scientific containers for mobility of compute. *PLoS One*. 2017;12(5):e0177459. <https://doi.org/10.1371/journal.pone.0177459>.
50. Wang J, Stephan R, Power K, Yan Q, Haehler H, Fanning S. Nucleotide sequences of 16 transmissible plasmids identified in nine multidrug-resistant *Escherichia coli* isolates expressing an ESBL phenotype isolated from food-producing animals and healthy humans. *J Antimicrob Chemother*. 2014;69(10):2658–68. <https://doi.org/10.1093/jac/dku206>.
51. Clement M, Ramette A, Bernasconi OJ, Principe L, Luzzaro F, Endimiani A. Whole-genome sequence of the first extended-spectrum beta-lactamase-producing strain of *Salmonella enterica subsp. enterica* serovar napoli. *Microbiol Resour Announc*. 2018;7(10):e00973.
52. Sydenham TV, Overballe-Petersen S, Hasman H, Wexler H, Kemp M, Justesen US. Complete hybrid genome assembly of clinical multidrug-resistant *Bacteroides fragilis* isolates enables comprehensive identification of antimicrobial-resistance genes and plasmids. *Microb Genom*. 2019;5(11): e000312.
53. Walker MJ, Brouwer S, Forde BM, Worthing KA, McIntyre L, Sundac L, et al. Detection of epidemic scarlet fever group A streptococcus in Australia. *Clin Infect Dis*. 2019;69(7):1232–4. <https://doi.org/10.1093/cid/ciz099>.
54. Taylor TL, Volkening JD, DeJesus E, Simmons M, Dimitrov KM, Tillman GE, et al. Rapid, multiplexed, whole genome and plasmid sequencing of foodborne pathogens using long-read nanopore technology. *Sci Rep*. 2019; 9(1):16350. <https://doi.org/10.1038/s41598-019-52424-x>.
55. Bessonov K, Robertson JA, Lin JT, Liu K, Gurnik S, Kernaghan SA, et al. Complete genome and plasmid sequences of 32 salmonella enterica strains from 30 serovars. *Microbiol Resour Announc*. 2018;7(17):e01232.
56. Pitt ME, Nguyen SH, Duarte TPS, Roddam LF, Blaskovich MAT, Cooper MA, et al. Complete genome sequences of clinical pandora fibrosis isolates. *Microbiol Resour Announc*. 2020;9(13):e00060.
57. Sieber RN, Overballe-Petersen S, Kaya H, Larsen AR, Petersen A. Complete genome sequences of methicillin-resistant staphylococcus aureus strains 110900 and 128254, two representatives of the CRISPR-cas-carrying sequence type 630/spa type t4549 lineage. *Microbiol Resour Announc*. 2020;9(41):e00891.
58. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat Biotechnol*. 2018;36(4):338–45. <https://doi.org/10.1038/nbt.4060>.
59. Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with multiplex MinION sequencing. *Microb Genom*. 2017;3(10): e000132. <https://doi.org/10.1099/mgen.0.000132>.
60. Wick RR, Holt KE. Benchmarking of long-read assemblers for prokaryote whole genome sequencing. *F1000Res*. 2019;8:2138.
61. Fu S, Wang A, Au KF. A comparative evaluation of hybrid error correction methods for error-prone long reads. *Genome Biol*. 2019;20(1):26. <https://doi.org/10.1186/s13059-018-1605-z>.
62. Wang L, Qu L, Yang L, Wang Y, Zhu H. NanoReviser: an error-correction tool for nanopore sequencing based on a deep learning algorithm. *Front Genet*. 2020;11:900. <https://doi.org/10.3389/fgene.2020.00900>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.