1    **Title: MicroPIPE:** *An end-to-end solution for high-quality complete bacterial genome*

2    *construction*

3

4    **Authors:** Valentine Murigneux[1]†, Leah W. Roberts[2,3,4]†, Brian M. Forde[5,6], Minh-Duy

5    Phan[5], Nguyen Thi Khanh Nhu[5], Adam D. Irwin[2,3], Patrick N. A. Harris[2,7], David L.

6    Paterson[2], Mark A. Schembri[5], David M. Whiley[2,3], Scott A. Beatson[5,6]

7

8    †These authors contributed equally

9

10   Affiliations:

11   1. QCIF Facility for Advanced Bioinformatics, Institute for Molecular Bioscience, The

12   University of Queensland, QLD, Australia

13   2. University of Queensland Centre for Clinical Research (UQCCR), Brisbane, Australia

14   3. Queensland Children's Hospital, Brisbane, Australia

15   4. European Bioinformatics Institute (EBI), European Molecular Biology Laboratory

16   (EMBL), Cambridge, Hinxton United Kingdom

17   5. School of Chemistry and Molecular Biosciences (SCMB), University of Queensland,

18   Brisbane, Australia

19   6. Australian Centre for Ecogenomics (ACE), University of Queensland, Brisbane, Australia

20   7. Central Microbiology, Pathology Queensland, Royal Brisbane & Women's Hospital,

21   Brisbane, Australia

22

23

24   Corresponding authors:

25   Leah Roberts

26   EMBL-EBI

27   leah@ebi.ac.uk

28

29   Scott Beatson

30   School of Chemistry and Molecular Biosciences, University of Queensland

31   scott.beatson@uq.edu.au

32

33

34

**Abstract:**

Oxford Nanopore Technology (ONT) long-read sequencing has become a popular platform for microbial researchers; however, easy and automated construction of high-quality bacterial genomes remains challenging. Here we present MicroPIPE: a reproducible end-to-end bacterial genome assembly pipeline for ONT and Illumina sequencing. To construct MicroPIPE, we evaluated the performance of several tools for genome reconstruction and assessed overall genome accuracy using ONT both natively and with Illumina. Further validation of MicroPIPE was carried out using 11 sequence type (ST)131 *Escherichia coli* and eight publicly available Gram-negative and Gram-positive bacterial isolates. MicroPIPE uses Singularity containers and the workflow manager Nextflow and is available at https://github.com/BeatsonLab-MicrobialGenomics/micropipe.

**Keywords:** Nanopore, ONT, pipeline, high quality, bacteria, assembly, polishing

**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 than through 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 (ONT), CCBGpipe [4], ASA³P [5] and Bactopia [6]. Katuali is an ONT developed assembly pipeline implemented in Snakemake. It offers the user flexibility in software choice, but with limited guidance or rationale. 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, this pipeline 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% [7] to the currently reported 95% for 1D reads using R9.4.1 flow cells [8]. While this is approaching that of Illumina (99.9%) [9] and PacBio (99%) [10], single nucleotide insertion/deletion (indel) errors remain problematic [11, 12]. 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 [13].
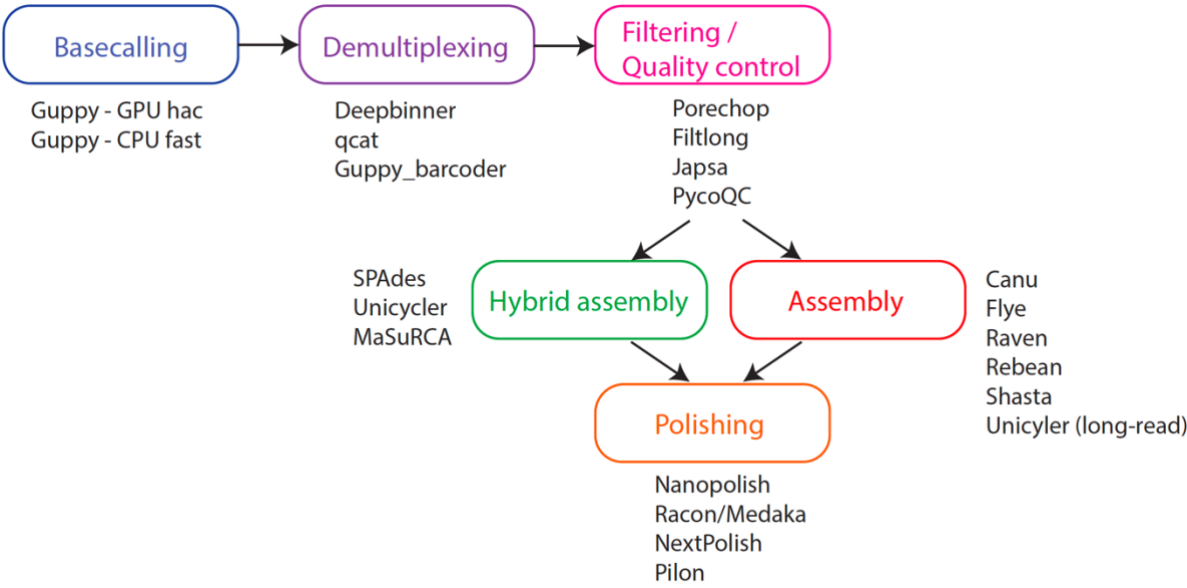
103   With the rapid pace of ONT progression, development of new software and pipelines, or
104   reappraisal of existing ones, has become an ongoing necessity. This has prompted the need for
105   appropriate validation sets, to assess (or reassess) the accuracy of results. While simulated
106   datasets provide an initial assessment of a tool's ability, data generated from biological sources
107   provide additional confidence in its real-world application, as has been developed previously
108   using metagenomic communities [14, 15]. *Escherichia coli* sequence type (ST)131 represents
109   a globally disseminated lineage that has been intensively studied as a result of its recent
110   emergence, antibiotic resistance and link to human disease [16-18]. Extensive knowledge of
111   both *E. coli* (as a species) and the ST131 lineage makes it an ideal dataset to use for software
112   and pipeline validation. Additionally, the *E. coli* ST131 strain EC958 represents an extensively
113   curated and highly accurate reference genome, having been sequenced on multiple occasions
114   using PacBio, Illumina and 454 pyrosequencing [19].

115

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

125

126   **Results:**

127

128   **Section 1: pipeline results and comparison to EC958 complete genome**

129

130   The main goal of this study was to create a robust and easily applicable pipeline for the
131   construction of high-quality bacterial genomes with minimal manual manipulations. To
132   achieve this, we first evaluated the performance of commonly used software at each stage of
133   bacterial genome construction using the high-quality EC958 genome (Accession: HG941718)
134   as our standard for final genome accuracy. **Figure 1** shows a diagram of the whole workflow,
135   indicating the software chosen for comparison at each stage. Nanopore reads for EC958 were
136   generated on a multiplexed run of 12 using the rapid barcoding kit on an R9.4.1 flow cell.

137



138

139 **Figure 1: overall diagram of assembly stages and tool comparisons**

140

141 ***Basecalling:***

142

143 To evaluate basecalling, we tested Guppy using both the "fast" and "high-accuracy" modes, as
144 well as the CPU vs. GPU configurations. When using Guppy v3.4.3 with the "high-accuracy"
145 setting on GPU servers we generated reads with approximately 91.0% accuracy in 828.5
146 minutes (13.81 hours). Using the "fast" mode on CPUs, we were able to generate 88.9%
147 accuracy in 2948.4 minutes (49.14 hours) **(Table 1)**. Testing the "high-accuracy" mode on a
148 CPU server was unfeasible due to the time required for processing (fewer than 10% of reads
149 completed basecalling in one week). Despite the lower per-read accuracy when using CPUs
150 and the "fast" basecalling setting, the consensus quality of the overall finished genome (after
151 assembly and polishing through MicroPIPE v0.8) was of comparable quality to that generated
152 with the GPU and high-accuracy setting **(Table 1)**.

153

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

158

159

160

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

| Basecalling comparison: | Guppy3.4.3_hac | Guppy3.4.3_fast | Guppy3.4.3_hac_ modbases | Guppy3.6.1_hac | Guppy3.6.1_hac_ modbases |
|---|---|---|---|---|---|
| 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 |

162

163

164    *Demultiplexing:*

165

166    For demultiplexing we tested three tools: Deepbinner, Guppy_barcoder and qcat. While Guppy

167    and qcat rely on basecalled reads, Deepbinner uses the raw fast5 reads. As such, we compared

168    the total numbers of binned reads after both basecalling and binning for each tool. Overall qcat

169    was the fastest demultiplexer, and was able to bin 89% of reads, compared to 84% for

170    Guppy_barcoder and 75% for Deepbinner **(Supplementary Figure 1)**. We prioritised read

171    retention to maximise coverage of each genome. As such, qcat was chosen as the default

172    demultiplexer for MicroPIPE. Following the recent depreciation of qcat, we have also provided

173    Guppy_barcoder as an optional demultiplexer.

174

175    *Filtering:*

176

177    Here we trialled two filtering tools: Filtlong and Japsa. Filtlong has the advantage of being

178    versatile enough to filter based on a number of requirements, such as read length, quality,

179    percentage of reads to keep and the option of using an external reference. Japsa primarily filters

180    based on read length and quality. Read metrics after filtering using each tool are given in

181    **Supplementary figure 2.** Overall, we found that filtering with Japsa retained more reads, but

182    with a reduced N50 read length and median read quality compared to Filtlong. Both tools took

183    an equivalent amount of time to run. For all downstream analysis we filtered reads using Japsa

184    with a minimum average quality cut-off of Q10 and 1 kb minimum read length, although

185    Filtlong would have been equally suitable. Both filtering tools are available as optional steps

186    in micropipe.

187

188    *Assembly:*

189

190    A number of tools have been designed for *de novo* assembly from long reads. Here we

191    compared six popular assembly tools and evaluated speed, completeness (of the chromosome

192    and plasmids, including circularisation) and correctness (i.e. nucleotide identity) based on the

193    complete EC958 reference genome standard, which contains 1 chromosome (5,109,767 bp)

194    and 2 plasmids (135,602 bp and 4080 bp). Parameters used for all assemblers are given in

195    **Supplementary Dataset 1**.

196

197    Overall, we found that all assemblers constructed the chromosome and larger (~135 kb)

198    plasmid **(Figure 2, Supplementary Table 2)**. Raven, Redbean and Shasta did not assemble

199    the smaller ~4 kb plasmid. While Canu was able to assemble both plasmids, closer inspection

200    found them to be much larger than expected (1.4x and 2x larger for the large and small plasmid,

201    respectively) due to overlapping ends that required additional trimming. Interestingly, both

202    Flye and Canu assembled a third, previously unidentified, small plasmid of ~1.8 kb in size.

203    This small plasmid was only identified when the Flye "—plasmids" mode was selected (to

204    rescue short unassembled plasmids) and when certain or no filtering parameters were applied

205    to the reads prior to assembly (**Supplementary Table 3**). Comparison of this small plasmid to
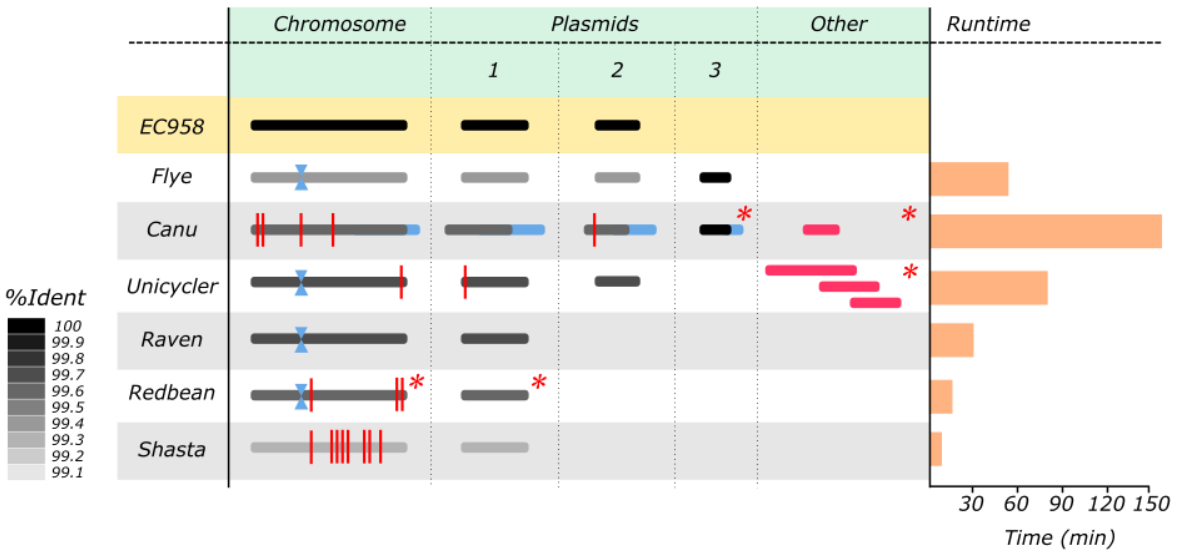
206   the Illumina data for the EC958 reference genome standard confirmed its presence and was

207   likely missed in the original assembly.

208

209   For most *de novo* assemblies, a number of small (<4.5 kb) misassemblies were detected, mainly

210   on the chromosome (**Figure 2**). This included a small inversion, which on closer inspection

211   was found to be an invertible phage tail protein that has been characterised previously [19].

212   This inversion was found in the Flye, Unicycler, Raven and Redbean assemblies and was not

213   counted as a misassembly due to its biological relevance.

214

215   Additional contigs were found in both Canu and Unicycler (long-read only mode). The three

216   additional contigs produced by Unicycler all matched other parts of the EC958 reference

217   genome standard (two on the chromosome, one on the larger plasmid). The additional contig

218   in Canu matched part of the additional ~1.8 kb plasmid.

219

220   In terms of speed, Shasta, Redbean and Raven were the fastest assemblers, completing in less

221   than 30 minutes. Of the remainder, Flye was four times faster than Canu and two times faster

222   than Unicycler. The majority of contigs from all assemblers were reported as circularised upon

223   assembly completion, with the exception of the additional contigs in Canu and Unicycler.

224   Redbean did not generate circularisation information, although the chromosome and plasmid

225   contigs could be circularised manually or using 3rd party software following assembly. Overall,

226   we found that Flye generated the best *de novo* assembly from long read data without the need
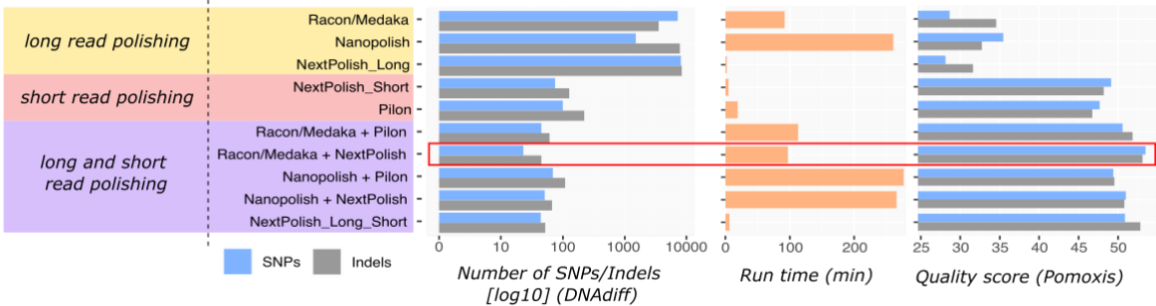
227   for manual intervention.

228

229

230

231

**Figure 2: Assembly comparison:** long horizonal bars represent contiguous sequences generated by each assembler. The chromosome and plasmids 1 and 2 are coloured according to their overall nucleotide identity when compared to the EC958 reference genome standard. The additional blue bars in the Canu plasmids represent the increased size of the plasmids from this assembler. Contigs that were not reported as circularised are marked with a red asterisk (*). Misassemblies are marked with a red vertical line at their approximate position. The phage tail protein inversion is marked with a blue hourglass.

*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 [13]. 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 (using long reads) followed by NextPolish (using short reads) achieved the most accurate assemblies **(Figure 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).

256

257  To confirm our choice of Flye as the best assembler, we polished assemblies generated from

258  the other five long-read assemblers, described above, using this strategy (**Supplementary**

259  **Table 5**). The polished Flye assembly remained the most accurate, closely followed by the
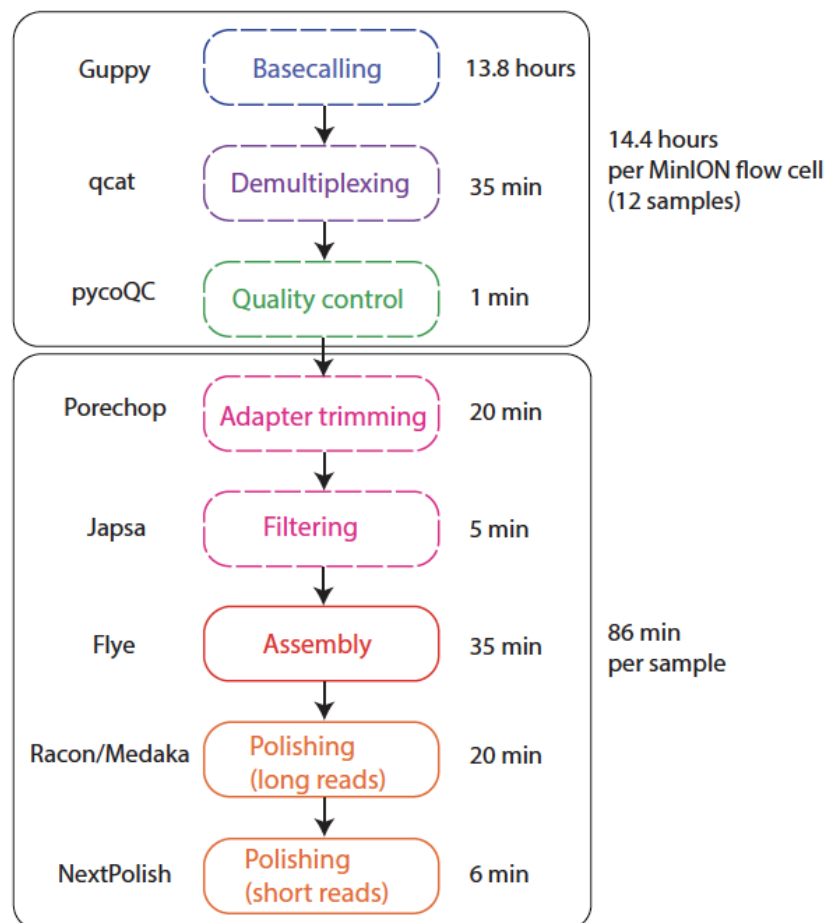
260  polished Raven assembly.

261



262

263  **Fig 3: Polishing results for EC958 ONT Flye assembly:** Comparative analysis of (i) long read polishing only,

264  (ii) short read polishing only, and (iii) sequential long read and short read polishing, using various tool

265  combinations. Comparison metrics were the number of SNPs/indels to the EC958 reference genome standard (by

266  DNAdiff), run time and quality score (by Poxomis assess_assembly).

267

268  *Hybrid assembly:*

269

270  In addition to long-read assembly (followed by short-read polishing), hybrid assemblers

271  capable of using both long and short reads simultaneously have also been developed, and

272  include Unicycler, MaSuRCA and SPAdes. Comparison of these pipelines to our genome

273  completed with Flye, Racon, Medaka and NextPolish found that they did not outperform our

274  current method. Unicycler was the only hybrid assembler able to completely resolve the

275  chromosome and both plasmids (SPAdes failed to circularise the chromosome while

276  MaSuRCA was unable to assemble the 4 kb plasmid) (**Supplementary Table 6**). Additional

277  long and short read polishing greatly improved the accuracy of the Unicycler and SPAdes

278  hybrid assemblies but not MaSuRCA (**Supplementary Table 5**). We compared the quality of

279  the genomes generated by either the best long-read only assembly (Flye) or the best hybrid

280  assembler based on accuracy and structure (Unicycler) and polished with the same strategy.

281  The polished assemblies contained a similar number of indels compared to the EC958 reference

282  genome standard, however the Flye assembly contained around two-fold fewer substitution

283  errors (**Supplementary Table 5**). Furthermore, Flye was nearly eight times faster than

284  Unicycler (**Supplementary Table 6**).

285

**Overall pipeline:**

287

288    Based on the results of our comparative analysis for all of the major steps of bacterial genome

289    assembly, we have developed MicroPIPE **(Figure 4)**. The pipeline is written in Nextflow [20]

290    and the dependencies are packaged into Singularity [21] container images available through

291    the Docker Hub and Quay.io BioContainers repositories. The bioinformatics workflow

292    manager Nextflow allows users to run the pipeline locally or using common High-Performance

293    Computing schedulers. Each step of the pipeline uses a specific container image which enables

294    easy modifications to be made in the future to include new or updated tools. The pipeline is

295    freely available on Github: https://github.com/BeatsonLab-MicrobialGenomics/micropipe.

296



297

298    **Figure 4: Overall pipeline**: Steps involved in genome assembly and the default tool selected for each stage.

299    Steps with dotted outline are optional. Time for running each step is provided based on running 12 multiplexed

300    *E. coli* samples with MicroPIPE v0.8. Basecalling (Guppy) and long-read polishing (Racon and Medaka) were

301    run on a GPU node. The rest of the pipeline was run using CPU resources.

302

**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 [19]. Overall, there were no other structural rearrangements. MicroPIPE assembled an additional ~1.8 kb plasmid, with 100% nucleotide identity to previously reported *E. coli* plasmids (GenBank records CP048320.1, KJ484633.1, [22]). 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, please 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 1-6 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**).

During preparation of this manuscript, Guppy v3.6.1 was released. MicroPIPE v0.9 (Guppy v3.6.1) was 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 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 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**).
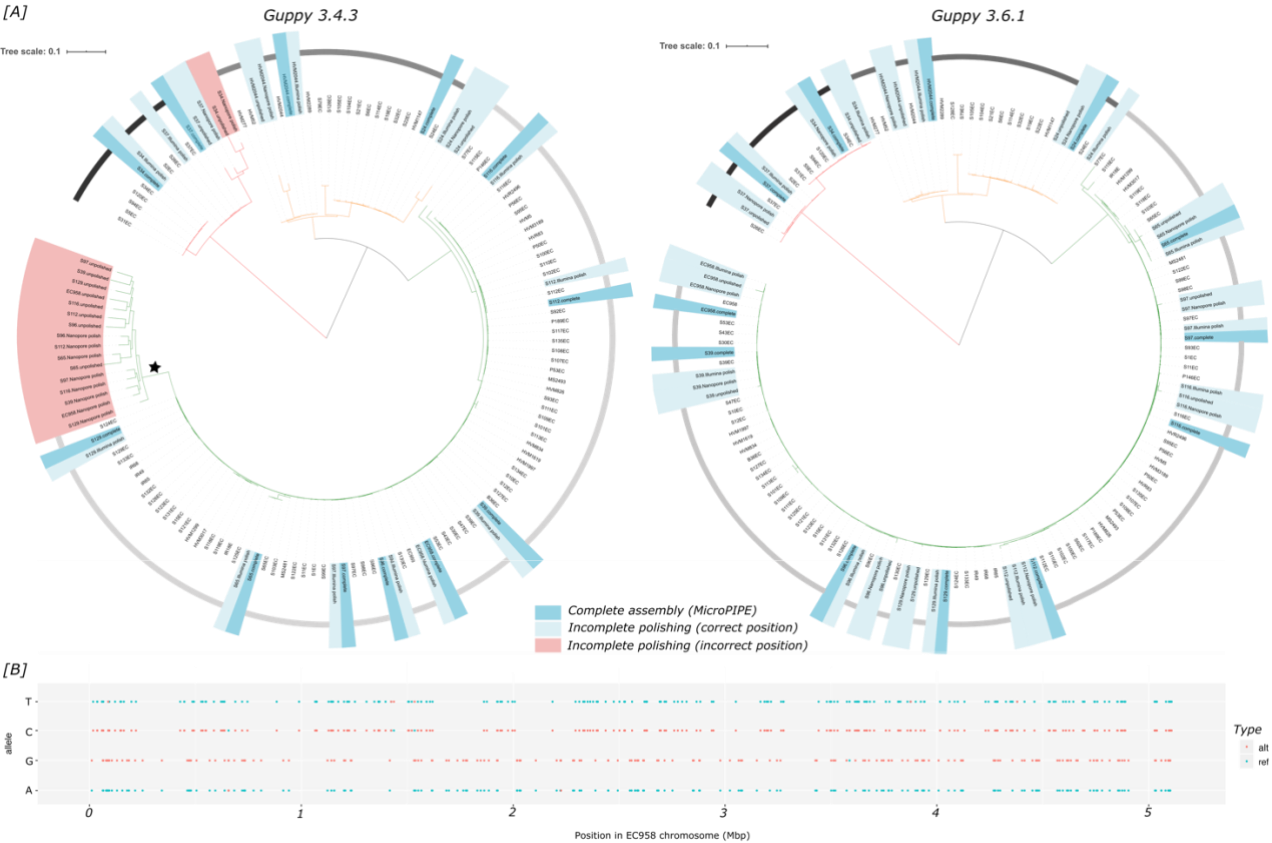
**Section 2: Validation of 11 ST131 *E. coli***

To further test the robustness of MicroPIPE on other genomes, we included an additional 11 well-characterised ST131 *E. coli* strains [16] on a multiplexed run of 12 *E. coli* (in addition to EC958).

Each strain took on average 86 minutes to run completely through MicroPIPE v0.8 using 16 threads (excluding the basecalling and demultiplexing steps) **(Figure 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 seven out of the 12 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 [16] 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 **(Figure 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 [16, 17] 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, Figure 5B)**. Further analysis of these sites determined that 393 (98%) were associated with a Dcm methylase motif CC(A/T)GG **(Supplementary Figure 4)**.

We also evaluated the MicroPIPE v0.9 assemblies using Guppy v3.6.1, which was released during preparation of this manuscript. By re-basecalling and recreating all assemblies as before,

370    we found a remarkable increase in the accuracy of Nanopore-only assemblies, such that all

371    assemblies clustered in their expected position within the tree **(Figure 5A)**.

372

373



374

375    **Figure 5: ST131 Phylogeny to assess quality of ONT assemblies: [A]** *dark blue*: Complete polished assemblies

376    from the MicroPIPE pipeline next to their Illumina assembly counterpart in the tree, *light blue*: assemblies with

377    incomplete polishing (i.e. Illumina only, Nanopore only or no polishing) clustered with their Illumina counterpart,

378    *red*: discrepant clustering of Nanopore assemblies. [B] position of alt alleles compared to the EC958 reference

379    standard chromosome present on branch leading to discrepant ONT assemblies as indicated by the star in (A).

380

381    **Section 3: MicroPIPE validation using publicly available ONT sequenced bacteria**

382

383    Lastly, we tested MicroPIPE using eight public genomes from both Gram-positive and Gram-

384    negative bacteria with available raw nanopore data (fast5) and validated our results using their

385    corresponding publicly available complete genomes (**Supplementary Dataset 1,**

386    **Supplementary Table 10**). As most of these isolates were sequenced using entire flow cells,

387    the coverage was reduced to 100x during the initial Flye assembly stage to minimise processing

388    time.

389

390   Using MicroPIPE v0.9, we were able to completely assemble the chromosome and plasmids

391   of all eight isolates. We were also able to recover two additional plasmids from the *Salmonella*

392   *enterica* str. SA20055162 that were not reported in the original assembly (**Supplementary**

393   **Table 10**).

394

395   To determine the accuracy of MicroPIPE, we compared our final assemblies with the submitted

396   complete genome for each isolate. Overall, the fewest differences were detected between our

397   MicroPIPE assembly and the complete genome of *Streptococcus pyogenes* strain SP1336,

398   constructed using PacBio long-read sequencing (8 SNPs, 96 indels; **Supplementary Table**

399   **10**). All other comparisons yielded 25-510 SNPs, and 14-758 indels, with the worst overall

400   being the *Salmonella enterica* serovar Napoli strain LC0541/17 (**Supplementary Table 10**).

401

402   With the exception of *S. pyogenes* SP1336, all other complete genomes were constructed using

403   previously assembled nanopore data (**Supplementary Dataset 1**). As such, we hypothesise

404   that our MicroPIPE assemblies likely represent corrections to the existing complete genomes,

405   as a result of updated basecalling and assembly methods. Further investigation found that one

406   sample, *Salmonella enterica* Bareilly str. CFSAN000189, also had a corresponding complete

407   genome constructed using PacBio data. Comparison of our MicroPIPE assembly to this

408   complete genome detected 0 SNPs and 15 indels, while there were 32 SNPs and 34 indels

409   compared to the ONT complete genome.

410

411   **Discussion:**

412

413   ONT long-read sequencing has quickly become one of the most prominent sequencing

414   platforms for microbial researchers globally. However, despite the large number of bacterial

415   genomes being completed using ONT, few end-to-end genome assembly pipelines exist. Here

416   we created an easy, automated and reproducible genome assembly pipeline for the construction

417   of complete, high-quality genomes using ONT in combination with Illumina sequencing. We

418   also provide a robust, publicly available set of 12 ST131 genomes that can be used to validate

419   future pipeline development or software advancements.

420

421   One of the main benefits of nanopore sequencing is its cost effectiveness, particularly when

422   multiplexing several samples onto a single flow cell. Methods have been developed to improve

423   yield and length during DNA extraction in order to achieve longer sequencing reads [14, 23].

424    However, here we show with our method that high-quality complete genomes can be achieved

425    using a standard, commercially available DNA extraction kit coupled with up to 12 multiplexed

426    samples. This build on other advances such as those described by Wick *et al.* [24], and

427    establishes an updated packaged pipeline that provides an efficient, cost effective and

428    reproducible approach to bacterial genome construction.

429

430    In our comparative analysis of different aspects of bacterial genome assembly, we chose not to

431    explore the effect of basecallers outside of ONT Guppy basecaller. This comparison has

432    already been completed previously [13], where it was found that Guppy outperformed other

433    existing basecallers. Guppy is also the default basecaller coupled with several of Oxford

434    Nanopore's devices, such as the MinIT, PromethION and GridION. For these reasons, we felt

435    that it was in the best interest of the community to provide a pipeline that used Guppy as the

436    basecaller. We also made a point of testing both the "high accuracy" mode on a GPU server

437    compared to the "fast" mode on a CPU server, as not all Nanopore users would have access to

438    GPU facilities. We found that, while the GPU server was significantly faster, basecalling reads

439    using the "fast" mode with CPUs can also achieve high-quality genomes with MicroPIPE.

440

441    During preparation of this manuscript, Guppy v3.6.1 was released with a raw read accuracy of

442    >97% using R9.4.1 flow cells (https://nanoporetech.com/accuracy). Community feedback

443    regarding this upgraded version supported increased overall accuracy, which prompted us to

444    incorporate this version into our analysis (MicroPipe v0.9). We also found that Guppy v3.6.1

445    increased the overall accuracy of our assemblies, particularly where it came to unresolved

446    indels using v3.4.3, which were suspected to be the result of technical artefacts around

447    methylated sites [23]. Using Guppy v3.6.1 made Nanopore-only assemblies more feasible,

448    particularly in cases where sufficient genetic context can be provided (e.g. identification of

449    outbreak vs. non-outbreak strains). However, we found that overall both v3.4.3 and v3.6.1 still

450    required polishing with short-read Illumina for maximum accuracy.

451

452    We observed some redundancy in the choice of tools for demultiplexing. Binning of reads with

453    both Guppy_barcoder and qcat performed almost equivalently (in terms of number of reads

454    binned), with minimal differences in the overall assembly **(Supplementary Table 11)**. Recent

455    improvements to Guppy_barcoder, which were released by ONT after compilation of this

456    manuscript, suggest that Guppy_barcoder is likely to be the default standard moving forward.

457

458    MicroPIPE implements a modest filtering measure to remove shorter, low quality reads from
459    the dataset. In this study, we found that filtering reads below 5 kb had little effect on the final
460    chromosome and larger plasmids, while filtering above 1 kb resulted in the loss of several small
461    plasmids in a number of strains (**Supplementary Tables 12 and 13**). Filtering with Filtlong at
462    "--min-length 1000 --keep_percent 90" resulted in the loss of the additional ~1.8 kb small
463    plasmid identified in EC958, which was retained when filtering with Japsa at "--min-length
464    1000". As such, we have implemented a 1 kb filtering cut-off (using Japsa) as default in
465    MicroPIPE to retain reads and small plasmids. However, we also found when testing
466    MicroPIPE on publicly available data that harsher filtering is sometimes desirable, especially
467    in cases where a single bacterial genome has been sequenced using an entire flow cell (such
468    that we used the Flye parameter "--asm-coverage 100" to reduce coverage for initial disjointig
469    assembly). As such, pre-processing of large quantities or highly ununiform data using Filtlong
470    may be the most desirable method. Ultimately, understanding the quality and read lengths of
471    the input data is a valuable step in generating the best possible assembly. We also provided the
472    user read quality assessment using PycoQC to assist in parameter selection.

473

474    Several other comparative analyses have been published exploring the overall utility of
475    different assemblers, in particular Wick *et al.* [25], who provide a comprehensive assembly
476    comparison using both simulated and real read datasets. While we did not test NECAT and
477    Miniasm, we found that our results generally matched those reported by Wick *et al.*,
478    particularly when it came to the overall strong performance of Flye. The most recent version
479    of Flye (v2.8) also removes the need to nominate a genome size, making it a more robust option.
480    However, we found that this version did not outperform the release used in this paper (v2.5)
481    on our dataset, as it was unable to circularise all plasmids. As such, we have retained Flye v2.5
482    in MicroPIPE.

483

484    Long and short read polishing is a staple of high-quality genome assembly, as the combination
485    of both ensures the correct contextual placement of variants as well as highly accurate
486    basecalls. However, while long-reads have enabled completion of assemblies by spanning
487    repetitive regions, polishing of these regions with short reads remains a problem. Here we
488    found that the majority of remaining differences between our EC958 ONT assembly and the
489    reference assembly (constructed with PacBio single molecule real time [SMRT] sequencing)
490    resided in repetitive regions. Ideally, polishing with long reads only would be a viable method
491    to reduce these errors as they would have sufficient coverage to ensure correct placement of

492    the repeat variant. However, as we show here, long read-only polishing was insufficient (likely

493    due to per-read accuracy), and short read polishing was necessary for removal of the majority

494    of errors. Currently, final polishing and assembly prior to completion will still necessitate

495    manual frameshift inspection. While impractical and costly, a combination of both PacBio and

496    ONT assembly could correct inherent biases in both technologies, using a consensus tool such

497    as Trycycler (https://github.com/rrwick/Trycycler). Long-read correction could also provide

498    another means of error reduction, however, was not assessed in this paper [26, 27].

499

500    We validated MicroPIPE using a set of 12 well-characterised *E. coli* isolates described

501    previously from a global collection [16, 17]. We did this for several reasons, including (i) the

502    availability of an existing high-quality reference genome and associated phylogenetic data (ii)

503    the robustness of *E. coli* as a representative species and workhorse organism, and (iii) our

504    extensive knowledge of the *E. coli* genome and ST131 lineage. We hope that by providing this

505    dataset to the wider community, it can serve as a resource for future validation and testing of

506    not only MicroPIPE, but other microbial assembly pipelines and tools.

507

508    In addition to in-house ONT sequencing data, we also tested MicroPIPE on a variety of publicly

509    available bacterial genomes to evaluate its assembly capabilities on other species. Without any

510    manual intervention, MicroPIPE was able to assemble all eight genomes, while also recovering

511    additional plasmids that were likely missed in the original assembly. When evaluating

512    correctness of the genomes, we found a number of remaining SNPs and indels when compared

513    to the complete genomes provided. Investigation into construction of the reference genomes

514    found that seven of the eight genomes provided were constructed previously using ONT

515    sequencing data, leading us to believe that differences in our assemblies compared to the

516    "reference" genomes may actually be corrections. Indeed, the genome with the closest match

517    between reference and MicroPIPE assembly were the genomes constructed using PacBio. As

518    such, we believe that genomes completed historically using ONT reads should be used

519    cautiously, and raw ONT data provided where possible to allow for reconstruction and

520    improvement of the assembly as the technology improves.

521

522    **Conclusions:**

523

524    Overall, we present an end-to-end pipeline for high-quality bacterial genome construction

525    designed to be easily implemented in the research lab setting. We believe this will be a useful

526  resource for users to easily and reproducibly construct complete bacterial genomes from

527  Nanopore sequencing data.

528

529  **Methods:**

530

531  **Public data:**

532  The EC958 complete genome was downloaded from NCBI (GenBank: HG941718.1,

533  HG941719.1, HG941720.1) [19]. Illumina reads for 12 ST131 genomes and draft assemblies

534  for 95 ST131 were accessed from [16]. Eight publicly available complete genomes were also

535  selected to test MicroPIPE, under the following criteria: (i) the raw nanopore sequencing files

536  (fast5) were available, (ii) a complete genome was made available for the same strain and (iii)

537  Illumina sequencing data were available for the same strain. These eight genomes represented

538  5 species from both gram-positive and gram-negative bacteria with chromosome sizes between

539  1.8 Mbp – 5.5 Mbps. A complete list of data used is provided in **Supplementary dataset 1**.

540

541  **Culture and DNA extraction:**

542  12 ST131 *E. coli* isolates (including EC958) were grown from single colonies in Lysogeny

543  Broth (LB) at 37ºC overnight with 250 rpm shaking. The overnight cultures (1.5 mL) were then

544  pelleted for DNA extraction using the Wizard Genomic DNA Purification Kit (Promega)

545  following manufacturer's protocol with modifications. Briefly, the cell pellet was lysed

546  following the protocol for Gram negative bacteria. RNA was removed by 1h incubation at 37°C

547  with RNase and the lysate was then mix with Protein Precipitation Solution by vortexing for

548  5s at max speed using Vortex-Genie 2 with horizontal tube adapter (Scientific Industries). The

549  DNA was precipitated using isopropanol and washed with 70% ethanol. The DNA pellet was

550  air-dried and then rehydrated in 100 µl EB buffer (QIAgen) by incubation at 65°C for 1 hour.

551  The DNA was quantified using a Qubit fluorometer (ThermoFisher Scientific) and the DNA

552  fragment size was estimated using agarose gel electrophoresis (0.5% agarose in TAE, 90V,

553  1h30m).

554

555  **Nanopore sequencing:**

556  DNA from 12 ST131 *E. coli* were multiplexed onto a single FLO-MIN106 flow cell using the

557  rapid barcode sequencing kit (SQK-RBK004) as per manufacturer's recommendation with the

558  following adjustments: the barcoded DNA was pooled without a concentration step using

559    AMPure XP beads prior to sequencing. Read metrics for each isolate are given in

560    **Supplementary table 1**.

561

562    **Pipeline tools and settings:**

563    Specific parameters and commands used to perform the following analyses are provided in full

564    in **Supplementary dataset 1**. MicroPIPE v0.8 uses Guppy v3.4.3, while MicroPIPE v0.9 uses

565    Guppy v3.6.1.

566

567    *Basecalling:*

568    Reads were basecalled using Guppy (v3.4.3) "fast" and "high-accuracy" modes. Fast mode was

569    evaluated using both GPU and CPU servers, while the "high-accuracy" mode was evaluated

570    using only GPU as the time to completion for this mode became unfeasible when run using

571    CPUs. Upon the release of Guppy v3.6.1, reads were re-basecalled using only the "high-

572    accuracy" mode. Guppy versions (3.4.3 and 3.6.1) were tested using the methylation aware

573    config file "dna_r9.4.1_450bps_modbases_dam-dcm-cpg_hac.cfg".

574

575    *Demultiplexing:*

576    Demultiplexing was evaluated using Guppy_barcoder (v3.4.3) and qcat (v1.0.1) on the

577    "passed" (>Q7) fastq reads after basecalling with Guppy. Demultiplexing using the raw fast5

578    reads was evaluated using Deepbinner (v0.2.0) [28]. Demultiplexed fast5 reads were

579    subsequently basecalled with Guppy (v3.4.3).

580

581    *Quality control:*

582    Barcodes and adapters were trimmed using Porechop (v0.2.3_seqan2.1.1)

583    (https://github.com/rrwick/Porechop). Overall read quality metrics and basecalling statistics

584    were extracted using PycoQC (v2.2.3) [29]. Read length and quality metrics per sample were

585    extracted using NanoPlot (v1.26.1) [30]. Average percentage read accuracy was determined by

586    mapping the basecalled reads to the reference genome EC958 using Minimap2 (v2.17-r954-

587    dirty) [31] and computing reads accuracy using Nanoplot. Filtering was evaluated using two

588    tools: Filtlong (v0.2.0) (https://github.com/rrwick/Filtlong) and Japsa (v1.9-01a)

589    (https://github.com/mdcao/japsa/).

590

591    *Assembly:*

592    Six assemblers were evaluated for long-read assembly only: Canu (v1.9) [32], Flye (v2.5) [33],

593    Raven (v1.1.5) (https://github.com/lbcb-sci/raven), Redbean (v2.5) [34], Shasta (v0.4.0: config

594    file                    optimised                    for                    Nanopore:

595    https://github.com/chanzuckerberg/shasta/blob/master/conf/Nanopore-Dec2019.conf)    [35]

596    and Unicycler (v0.4.7 long-read only) [36]. Three hybrid-assembly tools were also evaluated,

597    including SPAdes (v3.13.1) [37], Unicycler (v0.4.7) and MaSuRCA (v3.3.5) [38].

598

599    *Polishing and quality assessment*:

600    Polishing of the draft assemblies was evaluated using long reads (ONT), short reads (Illumina),

601    and a combination of both long and short reads. Long read polishing was performed using

602    Racon (v1.4.9) [39] and Medaka (v0.10.0) (https://nanoporetech.github.io/medaka/) (4

603    iterations of Racon based on Minimap2 v2.17-r941 overlaps followed by one iteration of

604    Medaka), Nanopolish (v0.11.1) [40] (1 iteration based on Minimap2 v2.17-r941 alignment)

605    and NextPolish (v1.1.0) [41] (2 iterations). Raw Illumina reads were trimmed using

606    Trimmomatic (v0.36) [42] with the following settings: ILLUMINACLIP:TruSeq3-PE-

607    2.fa:2:30:10 SLIDINGWINDOW:4:20 MINLEN:30. Short read polishing was performed

608    using NextPolish (v1.1.0) and Pilon (v1.23) [43] (both 2 rounds of polishing based on BWA

609    MEM v0.7.17-r1188 alignments).

610    Circularity was checked using NUCmer (v3.1) [44] to perform self-alignments. Final

611    assemblies were assessed for quality by comparison to the complete EC958 genome using the

612    assess_assembly tool from Pomoxis (v0.3) (https://github.com/nanoporetech/pomoxis) as well

613    as DNAdiff (v1.3) [44] and QUAST (v5.0.2) [45] to detect errors, misassemblies, and

614    determine overall nucleotide identity.

615

616    *Compute resources:*

617    All results were produced using cloud-based nodes with 16vCPUs and 32GB RAM. For the

618    GPU node, the GPU is a NVIDIA Tesla P40 24GB while the CPUs are 2x Intel Xeon Silver

619    4214 2.2G (12C/24T, 9.6GT/s, 16.5M Cache, Turbo, HT [85W] DDR4-2400).

620

621    **ST131 phylogeny:**

622    Parsnp (v1.5.2) [46] was used to create an ST131 phylogeny using the 12 ST131 *E. coli*

623    assembled in this study in addition to 95 ST131 *E. coli* short-read assemblies from Petty and

624    Ben Zakour *et al.* [16]. Recombination was removed using PhiPack [47], as implemented in

625    Parsnp. To evaluate the accuracy of each assembly and polishing step, we included our 12

626 completely polished assemblies (long and short read), 12 unpolished assemblies, 12 long-read

627 polished assemblies and 12 short-read polished assemblies. The tree was visualised using

628 Figtree (http://tree.bio.ed.ac.uk/software/figtree/) and iTOL [48].

629

630 **MEME methylation motif analysis:**

631 The 20 bps sequence (-10 to +10) around the 401 shared SNPs were extracted using BEDTools

632 getfasta (v2.28.0-33-g0f45761e) [49]. MEME (v5.2.0) [50, 51] was used to identify enriched

633 motifs within the sequences using the default parameters of the classic mode and allowing zero

634 or one occurrence per sequence. The motif CC(T/A)GG was significantly enriched in 393

635 sequences with an E-value of 6.2e-758.

636

637 **Declarations:**

638 **Ethics approval and consent to participate:**

639 Not applicable.

640

641 **Consent for publication:**

642 Not applicable.

643

644 **Availability of data and materials:**

645 The datasets generated and analysed during the current study are available under the following

646 Bioprojects (specific accessions available in supplementary dataset 1): EC958 complete

647 genome (GenBank: HG941718.1), ST131 Illumina data (PRJEB2968), ST131 Nanopore data

648 (fast5 and fastq [demultiplexed]; PRJNA679678).

649

650 **Competing interests:**

651 None to declare.

652

653 **Funding:**

654 LWR was supported by a Sakzewski Translational Research Grant. This work was supported

655 by funding from the Queensland Genomics Health Alliance (now Queensland Genomics),

656 Queensland Health, the Queensland Government.

657

658 **Authors' contributions:**

659 All authors conceptualised the study. VM, LWR, BMF and SAB developed the methodology.

660 MDP and MAS provided the bacterial strains and ONT sequencing data. VM wrote the

661 pipeline. VM, LWR and NTKN conducted formal analysis. All authors contributed to the

662 interpretation of results. SAB and MAS supervised aspects of the project and provided essential

663 expert analysis. LWR and VM wrote the original manuscript. BMF and SAB edited the

664 manuscript. All authors read and approved the final manuscript.

665

672

673 **References:**

674

675 1. Klassen JL, Currie CR: **Gene fragmentation in bacterial draft genomes: extent,**
676 **consequences and mitigation**. *BMC Genomics* 2012, **13**:14.
677 2. Koren S, Phillippy AM: **One chromosome, one contig: complete microbial genomes**
678 **from long-read sequencing and assembly**. *Curr Opin Microbiol* 2015, **23**:110-120.
679 3. Lemon JK, Khil PP, Frank KM, Dekker JP: **Rapid Nanopore Sequencing of Plasmids**
680 **and Resistance Gene Detection in Clinical Isolates**. *J Clin Microbiol* 2017,
681 **55**(12):3530-3543.
682 4. Liao YC, Cheng HW, Wu HC, Kuo SC, Lauderdale TY, Chen FJ: **Completing Circular**
683 **Bacterial Genomes With Assembly Complexity by Using a Sampling Strategy From a**
684 **Single MinION Run With Barcoding**. *Front Microbiol* 2019, **10**:2068.
685 5. Schwengers O, Hoek A, Fritzenwanker M, Falgenhauer L, Hain T, Chakraborty T,
686 Goesmann A: **ASA3P: An automatic and scalable pipeline for the assembly,**
687 **annotation and higher-level analysis of closely related bacterial isolates**. *PLoS*
688 *Comput Biol* 2020, **16**(3):e1007134.
689 6. Petit RA, 3rd, Read TD: **Bactopia: a Flexible Pipeline for Complete Analysis of**
690 **Bacterial Genomes**. *mSystems* 2020, **5**(4).
691 7. Rang FJ, Kloosterman WP, de Ridder J: **From squiggle to basepair: computational**
692 **approaches for improving nanopore sequencing read accuracy**. *Genome Biol* 2018,
693 **19**(1):90.
694 8. **R10.3: the newest nanopore for high accuracy nanopore sequencing – now**
695 **available in store** [https://nanoporetech.com/about-us/news/r103-newest-
696 nanopore-high-accuracy-nanopore-sequencing-now-available-store]
697 9. **Measuring sequencing accuracy**
698 [https://emea.illumina.com/science/technology/next-generation-sequencing/plan-
699 experiments/quality-scores.html]

700   10.   Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, Ebler J,
701         Fungtammasan A, Kolesnikov A, Olson ND *et al*: **Accurate circular consensus long-**
702         **read sequencing improves variant detection and assembly of a human genome**.
703         *Nat Biotechnol* 2019, **37**(10):1155-1162.
704   11.   Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q: **Opportunities and**
705         **challenges in long-read sequencing data analysis**. *Genome Biol* 2020, **21**(1):30.
706   12.   Weirather JL, de Cesare M, Wang Y, Piazza P, Sebastiano V, Wang XJ, Buck D, Au KF:
707         **Comprehensive comparison of Pacific Biosciences and Oxford Nanopore**
708         **Technologies and their applications to transcriptome analysis**. *F1000Res* 2017,
709         **6**:100.
710   13.   Wick RR, Judd LM, Holt KE: **Performance of neural network basecalling tools for**
711         **Oxford Nanopore sequencing**. *Genome Biol* 2019, **20**(1):129.
712   14.   Nicholls SM, Quick JC, Tang S, Loman NJ: **Ultra-deep, long-read nanopore**
713         **sequencing of mock microbial community standards**. *Gigascience* 2019, **8**(5).
714   15.   Sevim V, Lee J, Egan R, Clum A, Hundley H, Lee J, Everroad RC, Detweiler AM, Bebout
715         BM, Pett-Ridge J *et al*: **Shotgun metagenome data of a defined mock community**
716         **using Oxford Nanopore, PacBio and Illumina technologies**. *Sci Data* 2019, **6**(1):285.
717   16.   Petty NK, Ben Zakour NL, Stanton-Cook M, Skippington E, Totsika M, Forde BM, Phan
718         MD, Gomes Moriel D, Peters KM, Davies M *et al*: **Global dissemination of a**
719         **multidrug resistant Escherichia coli clone**. *Proc Natl Acad Sci U S A* 2014,
720         **111**(15):5694-5699.
721   17.   Ben Zakour NL, Alsheikh-Hussain AS, Ashcroft MM, Khanh Nhu NT, Roberts LW,
722         Stanton-Cook M, Schembri MA, Beatson SA: **Sequential Acquisition of Virulence and**
723         **Fluoroquinolone Resistance Has Shaped the Evolution of Escherichia coli ST131**.
724         *mBio* 2016, **7**(2):e00347-00316.
725   18.   Johnson JR, Porter S, Thuras P, Castanheira M: **The Pandemic H30 Subclone of**
726         **Sequence Type 131 (ST131) as the Leading Cause of Multidrug-Resistant**
727         **Escherichia coli Infections in the United States (2011-2012)**. *Open Forum Infect Dis*
728         2017, **4**(2):ofx089.
729   19.   Forde BM, Ben Zakour NL, Stanton-Cook M, Phan MD, Totsika M, Peters KM, Chan
730         KG, Schembri MA, Upton M, Beatson SA: **The complete genome sequence of**
731         **Escherichia coli EC958: a high quality reference sequence for the globally**
732         **disseminated multidrug resistant E. coli O25b:H4-ST131 clone**. *PLoS One* 2014,
733         **9**(8):e104400.
734   20.   Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C: **Nextflow**
735         **enables reproducible computational workflows**. *Nat Biotechnol* 2017, **35**(4):316-
736         319.
737   21.   Kurtzer GM, Sochat V, Bauer MW: **Singularity: Scientific containers for mobility of**
738         **compute**. *PLoS One* 2017, **12**(5):e0177459.
739   22.   Wang J, Stephan R, Power K, Yan Q, Hachler H, Fanning S: **Nucleotide sequences of**
740         **16 transmissible plasmids identified in nine multidrug-resistant Escherichia coli**
741         **isolates expressing an ESBL phenotype isolated from food-producing animals and**
742         **healthy humans**. *J Antimicrob Chemother* 2014, **69**(10):2658-2668.
743   23.   Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, Tyson JR, Beggs AD, Dilthey
744         AT, Fiddes IT *et al*: **Nanopore sequencing and assembly of a human genome with**
745         **ultra-long reads**. *Nat Biotechnol* 2018, **36**(4):338-345.

746   24.   Wick RR, Judd LM, Gorrie CL, Holt KE: **Completing bacterial genome assemblies with**
747         **multiplex MinION sequencing**. *Microb Genom* 2017, **3**(10):e000132.
748   25.   Wick RR, Holt KE: **Benchmarking of long-read assemblers for prokaryote whole**
749         **genome sequencing**. *F1000Res* 2019, **8**:2138.
750   26.   Fu S, Wang A, Au KF: **A comparative evaluation of hybrid error correction methods**
751         **for error-prone long reads**. *Genome Biol* 2019, **20**(1):26.
752   27.   Wang L, Qu L, Yang L, Wang Y, Zhu H: **NanoReviser: An Error-Correction Tool for**
753         **Nanopore Sequencing Based on a Deep Learning Algorithm**. *Front Genet* 2020,
754         **11**:900.
755   28.   Wick RR, Judd LM, Holt KE: **Deepbinner: Demultiplexing barcoded Oxford Nanopore**
756         **reads with deep convolutional neural networks**. *PLoS Comput Biol* 2018,
757         **14**(11):e1006583.
758   29.   Leger A, Leonardi T: **pycoQC, interactive quality control for Oxford Nanopore**
759         **Sequencing**. *Journal of Open Source Software* 2019, **4**(34).
760   30.   De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C: **NanoPack:**
761         **visualizing and processing long-read sequencing data**. *Bioinformatics* 2018,
762         **34**(15):2666-2669.
763   31.   Li H: **Minimap2: pairwise alignment for nucleotide sequences**. *Bioinformatics* 2018,
764         **34**(18):3094-3100.
765   32.   Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM: **Canu: scalable**
766         **and accurate long-read assembly via adaptive k-mer weighting and repeat**
767         **separation**. *Genome Res* 2017, **27**(5):722-736.
768   33.   Kolmogorov M, Yuan J, Lin Y, Pevzner PA: **Assembly of long, error-prone reads using**
769         **repeat graphs**. *Nat Biotechnol* 2019, **37**(5):540-546.
770   34.   Ruan J, Li H: **Fast and accurate long-read assembly with wtdbg2**. *Nat Methods* 2020,
771         **17**(2):155-158.
772   35.   Shafin K, Pesout T, Lorig-Roach R, Haukness M, Olsen HE, Bosworth C, Armstrong J,
773         Tigyi K, Maurer N, Koren S *et al*: **Nanopore sequencing and the Shasta toolkit enable**
774         **efficient de novo assembly of eleven human genomes**. *Nat Biotechnol* 2020,
775         **38**(9):1044-1053.
776   36.   Wick RR, Judd LM, Gorrie CL, Holt KE: **Unicycler: Resolving bacterial genome**
777         **assemblies from short and long sequencing reads**. *PLoS Comput Biol* 2017,
778         **13**(6):e1005595.
779   37.   Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,
780         Nikolenko SI, Pham S, Prjibelski AD *et al*: **SPAdes: a new genome assembly algorithm**
781         **and its applications to single-cell sequencing**. *J Comput Biol* 2012, **19**(5):455-477.
782   38.   Zimin AV, Puiu D, Luo MC, Zhu T, Koren S, Marcais G, Yorke JA, Dvorak J, Salzberg SL:
783         **Hybrid assembly of the large and highly repetitive genome of Aegilops tauschii, a**
784         **progenitor of bread wheat, with the MaSuRCA mega-reads algorithm**. *Genome Res*
785         2017, **27**(5):787-792.
786   39.   Vaser R, Sovic I, Nagarajan N, Sikic M: **Fast and accurate de novo genome assembly**
787         **from long uncorrected reads**. *Genome Res* 2017, **27**(5):737-746.
788   40.   Loman NJ, Quick J, Simpson JT: **A complete bacterial genome assembled de novo**
789         **using only nanopore sequencing data**. *Nat Methods* 2015, **12**(8):733-735.
790   41.   Hu J, Fan J, Sun Z, Liu S: **NextPolish: a fast and efficient genome polishing tool for**
791         **long-read assembly**. *Bioinformatics* 2020, **36**(7):2253-2255.

792    42.    Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**
793           **sequence data**. *Bioinformatics* 2014, **30**(15):2114-2120.
794    43.    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
795           Wortman J, Young SK *et al*: **Pilon: an integrated tool for comprehensive microbial**
796           **variant detection and genome assembly improvement**. *PLoS One* 2014,
797           **9**(11):e112963.
798    44.    Marcais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A: **MUMmer4: A**
799           **fast and versatile genome alignment system**. *PLoS Comput Biol* 2018,
800           **14**(1):e1005944.
801    45.    Gurevich A, Saveliev V, Vyahhi N, Tesler G: **QUAST: quality assessment tool for**
802           **genome assemblies**. *Bioinformatics* 2013, **29**(8):1072-1075.
803    46.    Treangen TJ, Ondov BD, Koren S, Phillippy AM: **The Harvest suite for rapid core-**
804           **genome alignment and visualization of thousands of intraspecific microbial**
805           **genomes**. *Genome Biol* 2014, **15**(11):524.
806    47.    Bruen TC, Philippe H, Bryant D: **A simple and robust statistical test for detecting the**
807           **presence of recombination**. *Genetics* 2006, **172**(4):2665-2681.
808    48.    Letunic I, Bork P: **Interactive Tree Of Life (iTOL) v4: recent updates and new**
809           **developments**. *Nucleic Acids Res* 2019, **47**(W1):W256-W259.
810    49.    Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic**
811           **features**. *Bioinformatics* 2010, **26**(6):841-842.
812    50.    Bailey TL: **Discovering novel sequence motifs with MEME**. *Curr Protoc*
813           *Bioinformatics* 2002, **Chapter 2**:Unit 2 4.
814    51.    Bailey TL, Elkan C: **Fitting a mixture model by expectation maximization to discover**
815           **motifs in biopolymers**. *Proc Int Conf Intell Syst Mol Biol* 1994, **2**:28-36.
816