# Comparison of long-read sequencing technologies in the hybrid assembly of complex bacterial genomes

Nicola De Maio[1]*, Liam P. Shaw[1]*, Alisdair Hubbard[2], Sophie George[1,3], Nick Sanderson[1], Jeremy Swann[1], Ryan Wick[4], Manal AbuOun[5], Emma Stubberfield[5], Sarah J. Hoosdally[1], Derrick W. Crook[1,3], Timothy E. A. Peto[1,3], Anna E. Sheppard[1,3], Mark J. Bailey[6], Daniel S. Read[6], Muna F. Anjum[5], A. Sarah Walker[1,3], Nicole Stoesser[1] on behalf of the REHAB consortium

1: Nuffield Department of Medicine, University of Oxford, Oxford, UK.
2: Department of Tropical Disease Biology, Liverpool School of Tropical Medicine, Liverpool, L3 5QA, UK.
3: HPRU IHR Health Protection Research Unit in Healthcare Associated Infections and Antimicrobial Resistance at University of Oxford in partnership with Public Health England, Oxford, UK.
4: Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, University of Melbourne, Australia.
5: Department of Bacteriology, Animal and Plant Health Agency, Addlestone, Surrey, KT15 3NB, UK.
6: Centre for Ecology & Hydrology, Benson Lane, Crowmarsh Gifford, Wallingford, OX10 8BB, UK.

* These authors contributed equally.

Corresponding authors: Nicole Stoesser, nicole.stoesser@ndm.ox.ac.uk

## ABSTRACT

26  
27  Illumina sequencing allows rapid, cheap and accurate whole genome bacterial analyses, but
28  short reads (<300 bp) do not usually enable complete genome assembly. Long read
29  sequencing greatly assists with resolving complex bacterial genomes, particularly when
30  combined with short-read Illumina data (hybrid assembly). However, it is not clear how
31  different long-read sequencing methods impact on assembly accuracy. Relative automation of
32  the assembly process is also crucial to facilitating high-throughput complete bacterial genome
33  reconstruction, avoiding multiple bespoke filtering and data manipulation steps. In this study,
34  we compared hybrid assemblies for 20 bacterial isolates, including two reference strains,
35  using Illumina sequencing and long reads from either Oxford Nanopore Technologies (ONT)
36  or from SMRT Pacific Biosciences (PacBio) sequencing platforms. We chose isolates from
37  the Enterobacteriaceae family, as these frequently have highly plastic, repetitive genetic
38  structures and complete genome reconstruction for these species is relevant for a precise
39  understanding of the epidemiology of antimicrobial resistance. We *de novo* assembled
40  genomes using the hybrid assembler Unicycler and compared different read processing
41  strategies. Both strategies facilitate high-quality genome reconstruction. Combining ONT and
42  Illumina reads fully resolved most genomes without additional manual steps, and at a lower
43  consumables cost per isolate in our setting. Automated hybrid assembly is a powerful tool for
44  complete and accurate bacterial genome assembly.

45

## IMPACT STATEMENT

47  Illumina short-read sequencing is frequently used for tasks in bacterial genomics, such as
48  assessing which species are present within samples, checking if specific genes of interest are
49  present within individual isolates, and reconstructing the evolutionary relationships between
50  strains. However, while short-read sequencing can reveal significant detail about the genomic
51  *content* of bacterial isolates, it is often insufficient for assessing genomic *structure*: how
52  different genes are arranged within genomes, and particularly which genes are on plasmids –
53  potentially highly mobile components of the genome frequently carrying antimicrobial
54  resistance elements. This is because Illumina short reads are typically too short to span
55  repetitive structures in the genome, making it impossible to accurately reconstruct these
56  repetitive regions. One solution is to complement Illumina short reads with long reads
57  generated with SMRT Pacific Biosciences (PacBio) or Oxford Nanopore Technologies (ONT)
58  sequencing platforms. Using this approach, called 'hybrid assembly', we show that we can
59  automatically fully reconstruct complex bacterial genomes of Enterobacteriaceae isolates in
60  the majority of cases (best-performing method: 17/20 isolates). In particular, by comparing
61  different methods we find that using the assembler Unicycler with Illumina and ONT reads
62  represents a low-cost, high-quality approach for reconstructing bacterial genomes using
63  publicly available software.

64

## DATA SUMMARY

66  Raw sequencing data and assemblies have been deposited in NCBI under BioProject
67  Accession PRJNA422511 (https://www.ncbi.nlm.nih.gov/bioproject/PRJNA422511). We
68  confirm all supporting data, code and protocols have been provided within the article or
69  through supplementary data files.

## INTRODUCTION

The rapid development of microbial genome sequencing methods over the last decade has revolutionized infectious disease epidemiology, and whole genome sequencing has become the standard for many molecular typing applications in research and public health (1–4). Much of this evolution has been driven by the development of high-throughput, low-cost, second generation (short-read) sequencing methods, such as Illumina's HiSeq and MiSeq platforms, which produce millions of low-error (0.1%) paired-end reads, generally 100-300bp in length. As such, Illumina sequencing has become the most widely used sequencing technology for microbial genomics. Multiple read processing algorithms now exist, typically enabling variant detection following mapping to a reference genome to assess genetic relatedness (e.g. for outbreak investigation or population genetic studies), or *de novo* assembly to facilitate the identification of important loci in the accessory genome, such as antimicrobial resistance genes (e.g. for epidemiological studies of resistance gene prevalence or for susceptibility prediction).

However, it has become clear that short-read sequencing has significant limitations depending on the bacterial species and/or epidemiological question. These limitations largely arise from the inability to fully reconstruct genomic structures of interest from short reads, including those on chromosomes and mobile genetic elements such as plasmids (5). An example where this genomic structure is highly relevant is the study of antimicrobial resistance (AMR) gene transmission and evolution in species of Enterobacteriaceae, which have emerged as a major clinical problem in the last decade (6). Short-read data from these species do not successfully facilitate assembly of the repetitive structures that extend beyond the maximum read length generated, including structures such as resistance gene cassettes, insertion sequences and transposons that are of crucial biological relevance to understanding the dissemination of key antimicrobial resistance genes.

The most widely used single molecule, long-read sequencing platforms, currently represented by Pacific Biosciences' (PacBio) Single Molecule Real-Time (SMRT) and Oxford Nanopore Technologies' (ONT) MinION sequencers, are often able to overcome these limitations by generating reads with a median length of 8-10kb, and as long as 100kb (5,7,8). However, the sequencing error rates of both long-read methods are much greater than Illumina (PacBio: 11-15%, raw, less in circular consensus reads (9); ONT: 5-40% (10)). Hybrid assembly, using combined short-read and long-read sequencing datasets, has emerged as a promising approach to generating fully resolved, accurate genome assemblies. With hybrid approaches, long reads provide information regarding the structure of the genome, specifically in plasmids, and short reads facilitate detailed assembly at local scales, and can be used to correct errors in long reads (11–13). The hybrid assembly tool Unicycler has been shown to outperform other hybrid assemblers in generating fully closed genomes (12).

We are not aware of any previously published direct comparisons of hybrid bacterial assemblies generated using long-read sequencing methods, yet the selection of a long-read sequencing approach has important cost, throughput and logistical implications. Currently, the two dominant long-read technologies are ONT and PacBio. The ONT MinION is a highly portable platform that has been deployed in several molecular laboratories, including those in low-income settings (14). Reported data yields of 10-30Gb and indexed barcoding now enable multiplexing of up to 12 bacterial isolates on a run (13). In contrast, the PacBio platform is non-portable but has been around longer, making it the most widely used for generating reference-grade bacterial assemblies to date (by way of example: as of 21st January 2019, NCBI Assembly contains 201 *E. coli* assemblies generated with PacBio vs. 3 generated with MinION).

121
122   Here we compared different approaches for hybrid bacterial genome assembly, using ONT
123   MinION, PacBio and Illumina HiSeq data generated from the same DNA extracts. We
124   selected 20 bacterial isolates from four genera of the Enterobacteriaceae family of bacteria
125   (*Escherichia*, *Klebsiella*, *Citrobacter* and *Enterobacter*) including two reference strains.
126   These genera typically have large bacterial genomes between 5-6.5Mb with diverse sets of
127   plasmids (15). We compared the advantages and disadvantages of ONT+Illumina versus
128   PacBio+Illumina hybrid assembly, including the need for additional manual processing steps.
129   We also investigated different strategies to optimize hybrid assembly using Unicycler for
130   both long-read approaches.

131

---

## METHODS

### Bacterial isolates, DNA extraction and Illumina sequencing

134   For sequencing, we selected and sub-cultured 20 isolates across the four genera of interest
135   from stocks of pure culture, stored in nutrient broth with 10% glycerol at        -80°C. Sub-
136   cultures were undertaken aerobically on Columbia blood agar at 37°C overnight. We chose
137   two reference strains, *Escherichia coli* CFT073, and *Klebsiella pneumoniae* MGH78578, and
138   18 isolates that were part of a study investigating antimicrobial resistance in diverse
139   Enterobacteriaceae from farm animals and environmental specimens (the REHAB study
140   http://modmedmicro.nsms.ox.ac.uk/rehab; details of isolates in Table S1). These comprised *E. coli*
141   (n=4), *K. pneumoniae* (n=2), *K. oxytoca* (n=2), *Citrobacter freundii* (n=2), *C. braakii* (n=2),
142   *C. gillenii* (n=1), *Enterobacter cloacae* (n=3), *E. kobei* (n=2). We chose to investigate
143   Enterobacteriaceae isolates as these bacteria are genetically complex: their genomes
144   commonly contain multiple plasmids and repeat structures of varying size, making them
145   difficult to assemble using other methods (5).

146

147   DNA was extracted from sub-cultured isolates using the Qiagen Genomic tip 100/G kit
148   (Qiagen, Valencia, CA, USA) to facilitate long-fragment extraction. Quality and fragment
149   length distributions were assessed using the Qubit fluorometer (ThermoFisher Scientific,
150   Waltham, MA, USA) and TapeStation (Agilent, Santa Clara, CA, USA).

151

152   All DNA extracts were sequenced using the Illumina HiSeq 4000, generating 150bp paired-
153   end reads. Libraries were constructed using the NEBNext Ultra DNA Sample Prep Master
154   Mix Kit (NEB, Ipswich, MA, USA) with minor modifications and a custom automated
155   protocol on a Biomek FX (Beckman Coulter, Brea, CA, USA). Ligation of adapters was
156   performed using Illumina Multiplex Adapters, and ligated libraries were size-selected using
157   Agencourt Ampure magnetic beads (Beckman Coulter, Brea, CA, USA). Each library was
158   PCR-enriched with custom primers (index primer plus dual index PCR primer (16)).
159   Enrichment and adapter extension of each preparation was obtained using 9μl of size-selected
160   library in a 50μl PCR reaction. Reactions were then purified with Agencourt Ampure XP
161   beads (Beckman Coulter, Brea, CA, USA) on a Biomek NXp after 10 cycles of amplification
162   (as per Illumina recommendations). Final size distributions of libraries were determined using
163   a TapeStation system as above and quantified by Qubit fluorometry.

164

### ONT library preparation and sequencing

166   ONT sequencing libraries were prepared by multiplexing DNA extracts from four isolates per
167   flowcell using the SQK-LSK108 and EXP-NBD103 kits according to the manufacturer's
168   protocol with the following amendments: input DNA (1.5μg) was not fragmented, 2ml
169   Eppendorf DNA LoBind tubes (Eppendorf, Hamburg, Germany) were used, all reactions
170   were purified using 0.4x Agencourt AMPure XP beads, incubation time with Agencourt

171 AMPure XP beads was doubled, elution volumes were reduced to the minimum required for
172 the subsequent step, and elution was heated to 37°C. Libraries were loaded onto flow cell
173 versions FLO-MIN106 R9.4 SpotON and sequenced for 48 hours.
174

## PacBio library preparation and sequencing

176 DNA extracts were initially sheared to an average length of 15kb using g-tubes, as specified
177 by the manufacturer (Covaris, Woburn, MA, USA). Sheared DNA was used in SMRTbell
178 library preparation, as recommended by the manufacturer. Quantity and quality of the
179 SMRTbell libraries were evaluated using the High Sensitivity dsDNA kit and Qubit
180 fluorometer and DNA 12000 kit on the 2100 Bioanalyzer (Agilent, Santa Clara, CA, USA).
181 To obtain the longest possible SMRTbell libraries for sequencing (as recommended by the
182 manufacturer), a further size selection step was performed using the PippinHT pulsed-field
183 gel electrophoresis system (Sage Science, Beverley, MA, USA), enriching for the SMRTbell
184 libraries >15kb for loading onto the instrument. Sequencing primer and P6 polymerase were
185 annealed and bound to the SMRTbell libraries, and each library was sequenced using a single
186 SMRT cell on the PacBio RSII sequencing system with 240-minute movies.
187

## Read preparation and hybrid assembly

189 ONT fast5 read files were base-called with Albacore (v2.0.2, https://github.com/JGI-
190 Bioinformatics/albacore), with barcode demultiplexing and fastq output. Adapter sequences were
191 trimmed with Porechop (v0.2.2, https://github.com/rrwick/Porechop). Read quality was calculated
192 with nanostat (v0.22, https://github.com/wdecoster/nanostat) (17).
193

194 Long reads from both ONT and PacBio were prepared using four alternative strategies:

195 • **Basic**: no filtering or correction of reads (i.e. all long reads available used for
196   assembly).

197 • **Corrected**: Long reads were error-corrected and subsampled (preferentially selecting
198   longest reads) to 30-40x coverage using Canu (v1.5, https://github.com/marbl/canu) (7)
199   with default options.

200 • **Filtered**: long reads were filtered using Filtlong (v0.1.1, https://github.com/rrwick/Filtlong)
201   by using Illumina reads as an external reference for read quality and either removing
202   10% of the worst reads or by retaining 500Mbp in total, whichever resulted in fewer
203   reads. We also removed reads shorter than 1kb and used the --trim and --split 250
204   options.

205 • **Subsampled**: we randomly subsampled long reads to leave approximately 600Mbp
206   (corresponding to a long read coverage around 100x).

207 Hybrid assembly for each of the two long-read sequencing technologies and for each of the
208 four read processing strategies (for a total of 8 hybrid assemblies per isolate) was performed
209 using Unicycler (v0.4.0) (12) with default options.
210

211 We used Bandage (v0.8.1) (18) to visualize assemblies, and the Interactive Genome Viewer
212 (IGV) (19) to visualize discrepancies in assemblies produced by the different methods.
213

214 To simulate the effect of additional multiplexing on ONT data and assembly (with current kits
215 allowing for up to 12 isolates to be indexed), we randomly subsampled half or one third of
216 the ONT reads from each isolate and repeated the assembly as in the "Basic" strategy above.
217

## Assembly comparison

219 We used multiple strategies to compare the features of different hybrid assemblies of the
220 same DNA extract. Firstly, we considered the completeness of an assembly i.e. specifically

221  whether all contigs reconstructed by Unicycler were identified as circular structures. Circular
222  structures typically represent completely assembled bacterial chromosomes and plasmids;
223  circular structures from different assemblies in our 20 isolates tended to agree in the majority
224  of cases (Table 1) and agreed with the structures of reference genomes for the two reference
225  strains (CFT073 and MH78578). We therefore also used the number of circular contigs in an
226  assembly as a measure of its completeness.

227

228  A common error associated with long-read-based assemblies is indel errors, which can
229  artificially shorten proteins by introducing premature stop codons or frameshift errors (20).
230  To check this possibility we annotated genomes with Prokka (v1.13.3,
231  https://github.com/tseemann/prokka) (21) then aligned all proteins to the full UniProt TrEMBL
232  database (November 15th 2018) using DIAMOND (v0.9.22, https://github.com/bbuchfink/diamond)
233  (22) and compared the length of each protein to its top hit. We compared proteins in
234  assemblies for the same sample with Roary (v3.12.0, https://sanger-pathogens.github.io/Roary) (23).

235

236  We additionally compared different assemblies of the same extract using:

237  • ALE (24), which assesses the quality of different assemblies using a likelihood-based
238    score of how well Illumina reads map to each assembly. ALE was run with default
239    parameters; Illumina reads were mapped to references using Bowtie2 (v2.3.3) (25).

240  • DNAdiff (as part of MUMMER v3.23) (26), which compares assemblies of the same
241    strain to detect differences such as SNPs and indels. DNAdiff was run with default
242    parameters on the fasta assembly files.

243  • REAPR (v1.0.18) (27), which (similarly to ALE) evaluates the accuracy of assemblies
244    using information from short read mapping to the assembly. REAPR was run using
245    the options "facheck", "smaltmap" and "pipeline" with default parameters.

246  • Minimap2 (v2017-09-21) (28) was used to map long reads to the hybrid assemblies,
247    and the mappings were evaluated to compare assembly quality and long read features
248    (identity and length) using scripts from the Filtlong package. We considered the
249    average identity for each base, and if there were multiple alignments at a base, we
250    used the one with the best score. We aligned PacBio and ONT reads to the hybrid
251    assemblies obtained either from all PacBio reads or from all ONT reads. Read
252    alignments were classified as: "good" if they had at least one alignment covering 97%
253    of the read, as a putative "chimera" if they had multiple inconsistent alignments
254    represented by at least 10% of the read length and ≥70% nucleotide identity, and
255    "other" if they did not fall into either of the two previous categories.

256

# RESULTS

## Sequencing data quality

259  For Illumina data, a median of 2,457,945 (interquartile range [IQR]: 2,073,342-2,662,727)
260  paired reads was generated for each isolate, with a median insert size of 363 bp (351-369).
261  The %GC content per isolate varied, as expected, by genus (median 53%, range: 50-57%),
262  but was consistent with the expected %GC content for each isolate based on its species (Table
263  S1).

264

265  The PacBio SMRT sequencing data resulted in a median of 160,740 (IQR: 153,196-169,240)
266  sub-reads with median sub-read length of 11,050 bp (IQR: 10,570-11,209 bp) per isolate.
267  Each isolate was sequenced using one SMRT cell on the RSII sequencing system, generating
268  a median of 1.32Gb (IQR: 1.25-1.36) of data per isolate, with isolates being run in batches of
269  8 (Figure S1, Table S1). For the ONT data, a median of 102,875 reads (IQR: 70,508-143,745

270  reads) were generated for each isolate, with a median phred score of 11.8 (IQR: 11.4-12.3).
271  ONT reads had a median length of 14,212 bp (IQR: 13,369-16,267 bp). A median of 13.8Gb
272  (IQR: 10.8-14.7Gb) of data was generated per run, resulting in a median of 3.45Gb per
273  isolate (four isolates multiplexed per run) (Figure S1, Table S1). After hybrid assembly, the
274  mean percentage identity and identity N50 for reads aligned against their respective
275  assemblies were higher for ONT reads than PacBio reads (mean±s.d. read alignment identity:
276  86±7 vs. 78±17; Figure S3, Table S3).
277
## Hybrid assembly runtimes
279  Clearly the computing infrastructure available to any given research team will be widely
280  variable, and assembly runtimes will therefore be different. For this experiment, where all
281  assemblies were run with dual 8-core Intel IvyBridge 2.6GHz, 256GB 1866MHz memory,
282  assembly times averaged between 1600-8000 minutes (~26-130 hours, Table S4), depending
283  on long-read preparation strategy (i.e. basic, corrected, filtered, sub-sampled, as in Methods).
284  They did not significantly vary depending on type of long-read used as input. Assemblies
285  completed in all cases, apart from two cases (both ONT+Illumina hybrids: MGH78578
286  reference strain, filtered strategy; RBHSTW-00123, corrected strategy).
287
## PacBio vs. ONT-based hybrid assembly comparisons
289  Using ONT+Illumina hybrid assembly approaches, we were able to completely assemble (i.e.
290  all contigs circularised) the majority of genomes (between 12 [60%] and 17 [85%] depending
291  on the preparation strategy for long reads, Table 1) without any manual intervention (18
292  across all strategies). With PacBio+Illumina fewer assemblies were complete (between 7
293  [35%] and 9 [45%]). More contigs were also circularised with ONT than with PacBio (up to
294  84 [97%] versus 67 [77%]), and assemblies were less fragmented (a minimum of 102 total
295  contigs across all 20 isolates for ONT vs. a minimum of 218 for PacBio).
296
297  On the basis of the minimap2/Filtlong comparisons (see Methods), most reads from both
298  long-read platforms had "good" alignment to their respective assemblies (~103,000 reads on
299  average for PacBio vs. ~99,000 reads for ONT, Figure S2, Table S2), with slightly more
300  alignments classified as "chimeras" (4,502 vs. 1,074 reads) and a much larger number of
301  alignments that were poor and classified as "other" (54,449 vs. 8,222) for PacBio compared
302  to ONT reads (Figure S2, Table S2).
303
304  Some chromosomal regions proved hard to assemble with both PacBio and ONT, e.g. for
305  isolates RBHSTW-00029 and RHB14-C01, but one of the noticeable differences between the
306  two methods was the ability of ONT to resolve repeats on small plasmids (see Figure 1 and
307  Figure S4). The DNA fragment size selection process used to optimize PacBio sequencing
308  and recommended by the manufacturer may have contributed to this (see Methods),
309  essentially making the assembly of small plasmids reliant on the Illumina short-read
310  component of the dataset only. This also reduces the power of PacBio reads for resolving the
311  genome structure when one copy of a repeated region is present on a short plasmid.
312
313  While correcting ONT reads with Canu or filtering them with Filtlong improved assembly
314  completeness for one isolate (RBHSTW-00309), in most cases avoiding this ONT read
315  correction and filtration led to better results (Table 1). This might be due to correction and
316  filtration steps removing reads in a non-uniform way across the genome, and in particular
317  from regions that are already hard to assemble. An alternative strategy deployed to reduce the
318  computational burden of hybrid assembly was to randomly sub-sample long reads until a
319  certain expected coverage was reached. Table 1 shows that this strategy was preferable to
320  read correction and filtration: it did not reduce assembly completeness but did reduce

321 computational demand (from an average of 5640 minutes to 2020 minutes per assembly on a
322 dual 8-core Intel IvyBridge 2.6GHz, 256GB 1866MHz memory, Table S4).
323
324 The analysis of local sequence assembly quality was inconclusive, showing inconsistent
325 results across different methodologies (Table 2), suggesting neither approach was clearly
326 superior to the other in this respect. However, detailed investigation of single nucleotide
327 polymorphisms (SNPs) between ONT and PacBio-based assemblies for the reference isolates
328 demonstrated two specific patterns of assembly differences. First, some positions (17 SNPs
329 across the two reference isolates) appeared plausibly polymorphic in the original DNA
330 sample and were called differently in different assembly runs (see Figure 2a). Secondly,
331 positions within regions with extremely low Illumina coverage (see Figure 2b) could have led
332 to assembly errors (25 SNPs across the two reference isolates), the PacBio assemblies being
333 more affected (22 cases vs 3 for ONT).
334
335 The proportion of proteins with a length of <90% of their top UniProt hit was low (~2-4% c.f.
336 3.7% for the RefSeq assembly of *E. coli* MG1655) and extremely consistent across
337 ONT+Illumina and PacBio+Illumina assemblies (Figure S5), suggesting that indels were not
338 a significant problem in the assemblies. There was very close agreement between methods
339 (median discrepancy < 5 proteins), although there were a greater number of cases where more
340 proteins were found in the ONT+Illumina assemblies (Figure S6). Proteins found uniquely in
341 an assembly tended to be found on a contig that was fragmented in the comparison assembly
342 (e.g. the third plasmid in the ONT-based assembly for RBHSTW-00167 was fragmented in
343 the comparison PacBio-based assembly, and was the location of 11 proteins unique to the
344 ONT-based assembly), highlighting that the degree of contig fragmentation in an assembly
345 can affect conclusions about gene presence beyond just the inability to resolve genomic
346 structures (Table S5, Figure S4).
347
348 Comparing *de novo* assemblies and reference genomes for the two reference strains (CFT073
349 and MGH78578) we found that the hybrid assemblies from ONT and PacBio reads were
350 more similar to each other (e.g. 18 SNPs and 0 indels for CFT073 and 24 SNPs and 13 indels
351 for MGH78578) than to the available reference genome sequences (156-365 SNPs and 47-
352 439 indels vs. the references, Table S6), possibly due to: (i) strain evolution in storage and
353 sub-culture since the reference strains were sequenced; (ii) errors in the original reference
354 sequences; and/or (iii) consistent errors in the hybrid assemblies.
355
356 Lastly, we investigated the effects of further ONT multiplexing by simulating datasets with 8
357 and 12 barcodes respectively (see Methods). Halving the available reads (equivalent to 8
358 barcodes) had no negative effect on the assemblies (Table S7). Using a third (equivalent to 12
359 barcodes) slightly increased the fragmentation of the assemblies overall (one fewer
360 completed assembly and nine additional non-circular contigs). However, these results were
361 not uniform: two assemblies gained an extra circular contig (RBHSTW-00309 and
362 RBHSTW-00340) with this downsampling.
363

## DNA preparation and sequencing costs

365 Beyond considerations of assembly accuracy, an important and realistic consideration when
366 choosing a sequencing approach is cost. While we do not attempt to calculate estimates that
367 will apply across different labs and settings, we can report our consumables costs per isolate
368 (i.e. exclusive of other potential costs, such as labour/infrastructure [laboratory and
369 computational]) in case it is helpful for informing others. The cost of bacterial culture and
370 DNA extraction was approximately £12 per isolate, resulting in sufficient DNA for all three
371 sequencing methods to be performed in parallel on a single extract. Cost for Illumina library
372 preparation and sequencing (see Methods) was ~£41 per isolate. ONT MinION sequencing

373 (library preparation and run) was performed by multiplexing 4 isolates per run, resulting in
374 costs of approximately £130 per isolate; however, it is possible to multiplex up to 12 isolates
375 per run at correspondingly lower coverage (13), resulting in costs of ~£44/isolate. At the time
376 we performed these experiments (late 2017), the PacBio sequencing was done using one
377 isolate per library per SMRTcell on the RSII system, with PacBio sequencing costs of more
378 than £280 per isolate. However, at the time of manuscript preparation, microbial sequencing
379 had been transferred to the higher throughput PacBio Sequel system, on which multiple
380 isolates can be multiplexed per SMRTcell 1M. Assuming ownership of a Sequel system, the
381 updated cost for PacBio sequencing, including DNA fragmentation, SMRTbell preparation,
382 size selection on the BluePippin system (Sage Science) and sequencing, is £190 per isolate
383 when multiplexing 8 isolates. If less coverage is needed or smaller genomes are to be
384 examined, one could multiplex up to 16 isolates per SMRTcell 1M at a cost of £152 per
385 isolate.

387 To summarise, in the optimal scenario for each technology in our setting, our total predicted
388 consumables costs range from £97-183 for generating an ONT+Illumina hybrid assembly
389 (multiplexing 4 versus 12 isolates) to £205-255 for generating a PacBio+Illumina hybrid
390 assembly on the PacBio Sequel system (multiplexing 8 versus 16 isolates). Costs using the
391 PacBio RSII system (i.e. >£320) to generate PacBio+Illumina hybrid assemblies would be
392 substantially higher than those for generating an ONT+Illumina hybrid assembly. We stress
393 that these costs are estimates only, and specifically do not include infrastructural and staffing
394 costs.

---

## DISCUSSION

397 Combining short read Illumina sequencing with different long read sequencing technologies
398 and using Unicycler, a publicly available and widely-used hybrid assembly tool, we found
399 that ONT+Illumina hybrid assembly generally facilitates the complete assembly of complex
400 bacterial genomes without additional manual steps. Our data thus support ONT+Illumina
401 sequencing as a non-inferior bacterial genome hybrid assembly approach compared with
402 PacBio+Illumina, leading to more complete assemblies, and to significantly lower costs per
403 isolate if multiplexed.

405 We also investigated the impact of different long-read processing strategies on assembly
406 quality and found that different strategies can result in more complete assemblies. We showed
407 that quality-based filtration and correction of long reads can apparently paradoxically result
408 in worse performance than just using unfiltered and uncorrected reads. There is no obvious
409 explanation for this; although we speculate that preferential removal of long reads from hard-
410 to-sequence regions might be a contributing factor, we have been unable to establish if this is
411 the case. We propose a different strategy to reduce the computational burden of hybrid
412 assembly without affecting the final outcome, namely randomly sub-sampling long reads
413 down to a desired level of coverage. We demonstrated that this strategy generally results in
414 better assemblies for ONT sequencing data.

416 We did however identify some recurrent patterns of local hybrid misassembly that could be
417 systematically addressed in the future. One of these is the presence of polymorphisms in the
418 DNA extract. These may represent genuine minor variants present in the isolate (although it
419 is difficult to establish with certainty), but the salient fact here is that current bacterial
420 assembly methods assume that no position is polymorphic which can lead to an imperfect
421 representation of the genomic content where this is not the case. We advocate for the
422 inclusion or awareness of polymorphisms within assembly polishing methods e.g. Pilon (29).
423

424　The other problem we identified is that regions with very low Illumina coverage tend to be
425　enriched with small assembly errors. This problem could similarly be addressed in the future
426　with hybrid assembly polishing methods, which would supplement Illumina-based polishing
427　with long read-based polishing in regions with low Illumina coverage.
428
429　There were several limitations to our study. Firstly, we included only two reference strains,
430　and our analyses suggest that the "true" sequences for these had diverged from the publicly
431　available reference sequences. This divergence could arise from multiple sources: true
432　biological variation after years of storage and/or sub-culture (a known possibility that has
433　been previously observed for bacterial reference strains e.g. in archived cultures of
434　*Salmonella enterica* serovar Typhimurium LT2 (30)), errors in the original reference
435　sequences (first published in 2002 for CFT073, 2007 for MGH78578), or possible errors in
436　our hybrid assemblies. Thus, making comparisons for any given approach even in the case
437　where a reference is available is difficult in the absence of a clear gold standard. Of note, we
438　tried to minimize biological variability introduced in culture by sequencing the same DNA
439　extract across different platforms. For 18 isolates the "true" underlying sequence was
440　unknown, which is common for highly plastic Enterobacteriaceae genomes. There is no
441　consensus on how best to evaluate assemblies and assembly quality when a reference is not
442　available. We therefore used several approaches, and these were not always consistent with
443　each other.
444
445　Assemblies can sometimes be further improved after an initial evaluation using "manual
446　completion" (see https://github.com/rrwick/Unicycler/wiki/Tips-for-finishing-genomes).
447　However, we did not investigate manual completion for our hybrid assemblies because, in
448　general, it is hard to replicate, has not been benchmarked and validated, is more easily biased,
449　and is not feasible for processing large numbers of isolates. We did not identify any
450　published, publicly available tools developed to specifically handle PacBio+Illumina hybrid
451　assembly, although some research groups may have implemented and validated these in-
452　house. Finally, we did not investigate the effect of different basecallers. The evolution of both
453　technologies and post-sequencing processing of data generated by both ONT and PacBio
454　platforms is rapid, and recent advances have been made e.g. in basecalling with the switch
455　from Albacore to Guppy for ONT data. Our assumption is that such advances which improve
456　read quality and basecalling will improve assembly quality, but we have not carried out
457　specific comparisons.
458
459　In conclusion, we have demonstrated that reference-grade, complete hybrid assemblies can be
460　effectively generated for complex bacterial genomes including multiple plasmids using ONT
461　platforms in combination with Illumina data. Given the average yields that can be generated
462　with these devices, it should be feasible to comfortably multiplex eight Enterobacteriaceae
463　isolates per ONT flowcell. At current listed cost prices, this effectively represents a cost of
464　~£100/hybrid assembly (all laboratory and sequencing consumables costs [includes Illumina
465　and Nanopore]).
466
467

## AUTHOR STATEMENTS

469
### Authors and contributors
471　Conceptualisation: NdM, ASW, TEAP, DWC, NSt; Methodology: NdM, LPS, RW, NSt;
472　Software: NdM, LPS, RW, AS, NSa, JS; Formal analysis: NdM, LPS, NSt; Investigation: AH,
473　SG, MAb, ES; Resources: MA, DR, DWC, ASW, TEAP, SJH, NSt; Data curation: MAb, ES,
474　NdM, LPS, NSt; Writing - original draft preparation: NdM, LPS, SG, ASW, NSt; Writing -

475 review and editing: All authors; Visualisation: NdM, LPS, NSt; Supervision: MA, DR, DWC,
476 NSt; Project administration: SJH, MA, DR, NSt; Funding: MA, DR, DWC, NSt.
477

## ABBREVIATIONS

511

513 ONT: Oxford Nanopore Technologies
514 PacBio: Pacific Biosciences
515 SNP: single nucleotide polymorphism
516 AMR: antimicrobial resistance
517

## REFERENCES

519

520 1. Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. Transforming clinical
521 microbiology with bacterial genome sequencing. *Nature Reviews. Genetics*. [Online]
522 2012;13(9): 601–612. Available from: doi:10.1038/nrg3226

523   2.   Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, et al. Rapid antibiotic-
524        resistance predictions from genome sequence data for Staphylococcus aureus and
525        Mycobacterium tuberculosis. *Nature Communications*. [Online] 2015;6: 10063.
526        Available from: doi:10.1038/ncomms10063

527   3.   Didelot X, Walker AS, Peto TE, Crook DW, Wilson DJ. Within-host evolution of
528        bacterial pathogens. *Nature Reviews. Microbiology*. [Online] 2016;14(3): 150–162.
529        Available from: doi:10.1038/nrmicro.2015.13

530   4.   Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut
531        microbial gene catalogue established by metagenomic sequencing. *Nature*. [Online]
532        2010;464(7285): 59–65. Available from: doi:10.1038/nature08821

533   5.   George S, Pankhurst L, Hubbard A, Votintseva A, Stoesser N, Sheppard AE, et al.
534        Resolving plasmid structures in Enterobacteriaceae using the MinION nanopore
535        sequencer: assessment of MinION and MinION/Illumina hybrid data assembly
536        approaches. *Microbial Genomics*. [Online] 2017;3(8): e000118. Available from:
537        doi:10.1099/mgen.0.000118

538   6.   Logan LK, Weinstein RA. The Epidemiology of Carbapenem-Resistant
539        Enterobacteriaceae: The Impact and Evolution of a Global Menace. *The Journal of
540        Infectious Diseases*. [Online] 2017;215(suppl_1): S28–S36. Available from:
541        doi:10.1093/infdis/jiw282

542   7.   Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable
543        and accurate long-read assembly via adaptive k-mer weighting and repeat separation.
544        *Genome Research*. [Online] 2017;27(5): 722–736. Available from:
545        doi:10.1101/gr.215087.116

546   8.   Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using
547        only nanopore sequencing data. *Nature Methods*. [Online] 2015;12(8): 733–735.
548        Available from: doi:10.1038/nmeth.3444

549   9.   Rhoads A, Au KF. PacBio Sequencing and Its Applications. *Genomics, Proteomics &
550        Bioinformatics*. [Online] 2015;13(5): 278–289. Available from:
551        doi:10.1016/j.gpb.2015.08.002

552   10.  Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz MC, McCombie WR.
553        Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a
554        eukaryotic genome. *Genome Research*. [Online] 2015;25(11): 1750–1756. Available
555        from: doi:10.1101/gr.191395.115

556   11.  Risse J, Thomson M, Patrick S, Blakely G, Koutsovoulos G, Blaxter M, et al. A single
557        chromosome assembly of Bacteroides fragilis strain BE1 from Illumina and MinION
558        nanopore sequencing data. *GigaScience*. [Online] 2015;4: 60. Available from:
559        doi:10.1186/s13742-015-0101-6

560   12.  Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome
561        assemblies from short and long sequencing reads. *PLoS computational biology*. [Online]
562        2017;13(6): e1005595. Available from: doi:10.1371/journal.pcbi.1005595

563　13. Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with
564　　　multiplex MinION sequencing. *Microbial Genomics*. [Online] 2017;3(10): e000132.
565　　　Available from: doi:10.1099/mgen.0.000132

566　14. Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, et al. Real-time,
567　　　portable genome sequencing for Ebola surveillance. *Nature*. [Online] 2016;530(7589):
568　　　228–232. Available from: doi:10.1038/nature16996

569　15. Carattoli A. Resistance plasmid families in Enterobacteriaceae. *Antimicrobial Agents*
570　　　*and Chemotherapy*. [Online] 2009;53(6): 2227–2238. Available from:
571　　　doi:10.1128/AAC.01707-08

572　16. Lamble S, Batty E, Attar M, Buck D, Bowden R, Lunter G, et al. Improved workflows
573　　　for high throughput library preparation using the transposome-based Nextera system.
574　　　*BMC biotechnology*. [Online] 2013;13: 104. Available from: doi:10.1186/1472-6750-13-
575　　　104

576　17. De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack:
577　　　visualizing and processing long-read sequencing data. *Bioinformatics*. [Online]
578　　　2018;34(15): 2666–2669. Available from: doi:10.1093/bioinformatics/bty149

579　18. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo
580　　　genome assemblies. *Bioinformatics (Oxford, England)*. [Online] 2015;31(20): 3350–
581　　　3352. Available from: doi:10.1093/bioinformatics/btv383

582　19. Thorvaldsdottir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-
583　　　performance genomics data visualization and exploration. *Briefings in Bioinformatics*.
584　　　[Online] 2013;14(2): 178–192. Available from: doi:10.1093/bib/bbs017

585　20. Watson M. A simple test for uncorrected insertions and deletions (indels) in bacterial
586　　　genomes. 2018; Available from: http://www.opiniomics.org/a-simple-test-for-
587　　　uncorrected-insertions-and-deletions-indels-in-bacterial-genomes/

588　21. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics (Oxford,*
589　　　*England)*. [Online] 2014;30(14): 2068–9. Available from:
590　　　doi:10.1093/bioinformatics/btu153

591　22. *Fast and sensitive protein alignment using DIAMOND | Nature Methods*. [Online]
592　　　Available from: https://www.nature.com/articles/nmeth.3176 [Accessed: 24th January
593　　　2019]

594　23. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary: Rapid
595　　　large-scale prokaryote pan genome analysis. *Bioinformatics*. [Online] 2015;31(22):
596　　　3691–3693. Available from: doi:10.1093/bioinformatics/btv421

597　24. Clark SC, Egan R, Frazier PI, Wang Z. ALE: a generic assembly likelihood evaluation
598　　　framework for assessing the accuracy of genome and metagenome assemblies.
599　　　*Bioinformatics (Oxford, England)*. [Online] 2013;29(4): 435–443. Available from:
600　　　doi:10.1093/bioinformatics/bts723

601　25. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods*.
602　　　[Online] 2012;9(4): 357–359. Available from: doi:10.1038/nmeth.1923

603   26.   Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile
604         and open software for comparing large genomes. *Genome Biology*. [Online] 2004;5(2):
605         R12. Available from: doi:10.1186/gb-2004-5-2-r12

606   27.   Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD. REAPR: a universal
607         tool for genome assembly evaluation. *Genome Biology*. [Online] 2013;14(5): R47.
608         Available from: doi:10.1186/gb-2013-14-5-r47

609   28.   Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. [Online]
610         2018;34(18): 3094–3100. Available from: doi:10.1093/bioinformatics/bty191

611   29.   Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an
612         integrated tool for comprehensive microbial variant detection and genome assembly
613         improvement. *PloS One*. [Online] 2014;9(11): e112963. Available from:
614         doi:10.1371/journal.pone.0112963

615   30.   Rabsch W, Helm RA, Eisenstark A. Diversity of phage types among archived cultures of
616         the Demerec collection of Salmonella enterica serovar Typhimurium strains. *Applied
617         and Environmental Microbiology*. 2004;70(2): 664–669.

618
619

# FIGURES AND TABLES

**Table 1. Summary of all assemblies in terms of circularised contigs.** Different rows refer to different isolates. "*n* of *m*" means that *n* contigs were circular in the assembly out of *m* total contigs. When *n* and *m* are identical, it means that the assembly was considered complete, and these cases are shaded in green. "Basic", "Corrected", "Filtered" and "Subsampled" refer to the strategies of long read preparation (see Methods). "NA" refers to cases where the assembly pipeline repeatedly failed. The true number of circular structures was estimated by inspection.
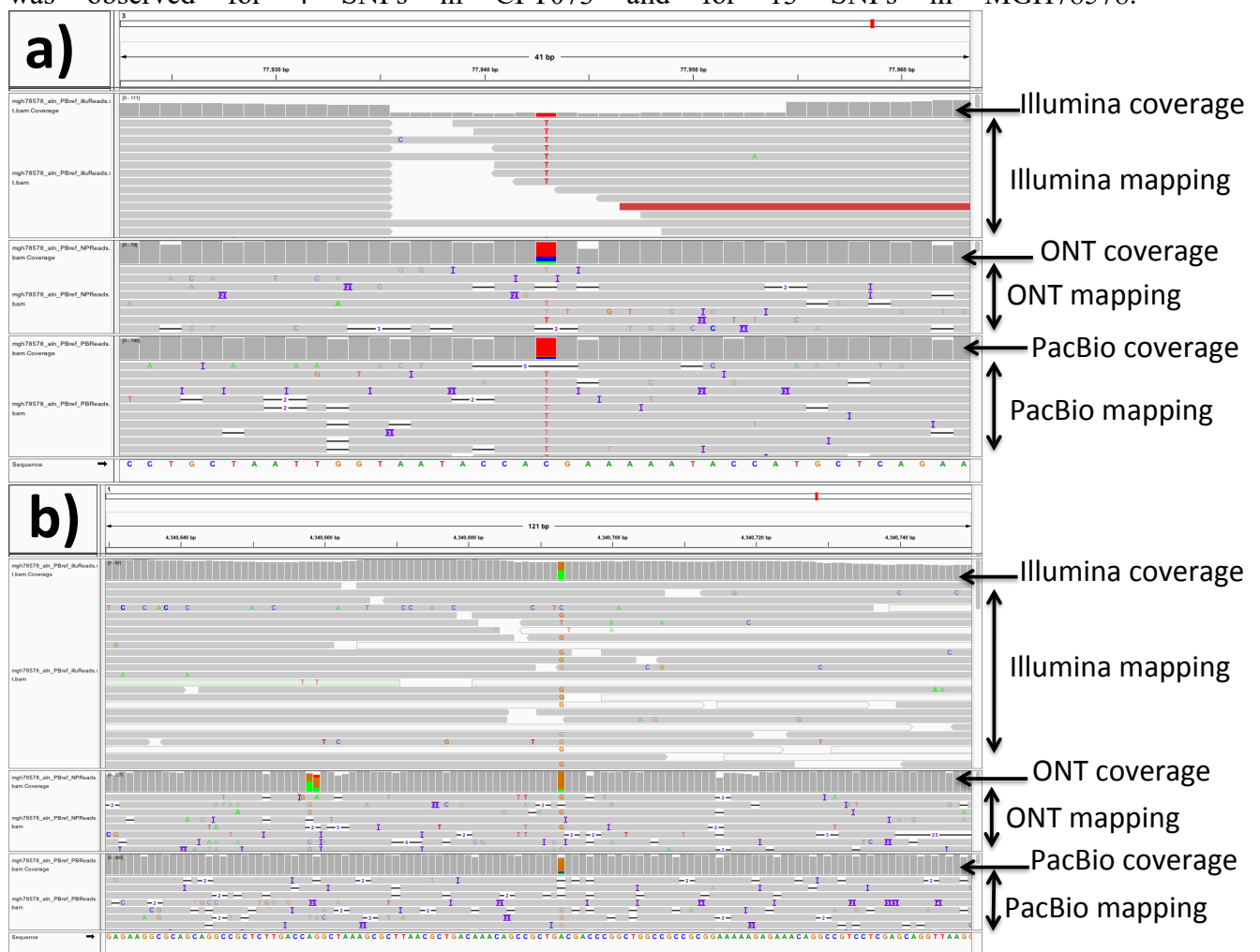
| Isolate | ONT (MinION) | | | | PacBio (RSII System) | | | | True circular structures (estimated) |
|---|---|---|---|---|---|---|---|---|---|
| | Basic | Corrected | Filtered | Subsampled | Basic | Corrected | Filtered | Subsampled | |
| CFT073 (reference) | 1 of 1 | 1 of 1 | 0 of 9 | 1 of 1 | 0 of 9 | 0 of 9 | 0 of 9 | 0 of 9 | 1 |
| MGH78578 (reference) | 6 of 6 | 4 of 7 | NA | 6 of 6 | 4 of 7 | 2 of 22 | 2 of 22 | 2 of 22 | 6 |
| RBHSTW-00029 | 3 of 9 | 3 of 9 | 3 of 9 | 3 of 9 | 3 of 9 | 3 of 9 | 3 of 9 | 3 of 9 | 4 |
| RBHSTW-00053 | 6 of 6 | 6 of 6 | 6 of 6 | 6 of 6 | 6 of 6 | 6 of 6 | 6 of 6 | 6 of 6 | 6 |
| RBHSTW-00059 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 |
| RBHSTW-00122 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 |
| RBHSTW-00123 | 7 of 7 | NA | 7 of 7 | 7 of 7 | 5 of 8 | 4 of 18 | 4 of 18 | 4 of 18 | 7 |
| RBHSTW-00127 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 of 5 | 5 |
| RBHSTW-00128 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 4 of 4 | 3 of 6 | 3 of 6 | 3 of 6 | 4 |
| RBHSTW-00131 | 4 of 4 | 2 of 7 | 4 of 4 | 4 of 4 | 3 of 15 | 4 of 5 | 3 of 15 | 2 of 15 | 4 |
| RBHSTW-00142 | 7 of 7 | 5 of 25 | 7 of 7 | 7 of 7 | 4 of 24 | 4 of 58 | 4 of 24 | 4 of 27 | 7 |
| RBHSTW-00167 | 9 of 9 | 5 of 15 | 10 of 10 | 9 of 9 | 4 of 34 | 3 of 60 | 3 of 60 | 3 of 60 | 9 |
| RBHSTW-00189 | 6 of 6 | 6 of 6 | 5 of 6 | 6 of 6 | 5 of 29 | 5 of 28 | 5 of 29 | 5 of 30 | 6 |
| RBHSTW-00277 | 2 of 2 | 2 of 2 | 1 of 8 | 2 of 2 | 1 of 8 | 1 of 8 | 1 of 8 | 1 of 8 | 2 |
| RBHSTW-00309 | 4 of 5 | 5 of 5 | 5 of 5 | 4 of 5 | 5 of 5 | 4 of 5 | 5 of 5 | 5 of 5 | 5 |
| RBHSTW-00340 | 3 of 11 | 3 of 11 | 4 of 4 | 4 of 4 | 2 of 25 | 2 of 25 | 2 of 24 | 2 of 25 | 4 |
| RBHSTW-00350 | 2 of 2 | 2 of 2 | 2 of 3 | 2 of 2 | 2 of 2 | 2 of 2 | 2 of 2 | 2 of 2 | 2 |
| RHB10-C07 | 1 of 1 | 1 of 1 | 1 of 1 | 1 of 1 | 1 of 1 | 1 of 1 | 1 of 17 | 1 of 1 | 1 |
| RHB11-C04 | 3 of 3 | 3 of 3 | 3 of 3 | 3 of 3 | 3 of 3 | 3 of 3 | 3 of 3 | 3 of 3 | 3 |
| RHB14-C01 | 1 of 12 | 1 of 12 | 1 of 15 | 1 of 12 | 1 of 15 | 1 of 15 | 1 of 15 | 1 of 15 | 2 |
| Total contigs | 109 | 130 | 115 | 102 | 218 | 294 | 276 | 265 | |

| | ONT (MinION) | | | | PacBio (RSII System) | | | |
|---|---|---|---|---|---|---|---|---|
| Total circularised contigs (% over total estimated circular structures from Bandage: n=87 for all isolates) | 83 (95%) | 67 (84%) | 77 (95%) | 84 (97%) | 67 (77%) | 62 (71%) | 62 (71%) | 61 (70%) |
| Total circularised contigs for reference strains (i.e. structures known, total n=1 [*E. coli*] + 6 [*K. pneumoniae*]) | 7 (100%) | 5 (71%) | 0 (0%) | 7 (100%) | 5 (71%) | 2 (29%) | 2 (29%) | 2 (29%) |
| Total isolates with all contigs circularised (% isolates) | 16 (80%) | 12 (60%) | 13 (65%) | 17 (85%) | 9 (45%) | 7 (35%) | 7 (35%) | 8 (40%) |

625 **Table 2. Comparison between PacBio and ONT-based hybrid assemblies**. Comparisons are shown using ALE, DNAdiff and REAPR (see
626 Methods). Different rows represent different isolates. All entries representing a better score for the PacBio assembly are shaded in red, those
627 showing a better score for ONT are shaded in blue. "ALE score" is the assembly likelihood difference (calculated by ALE from the mapping of
628 Illumina reads) between PacBio and ONT assemblies. "Unmapped reads" refers to the number of Illumina reads that ALE did not map to the
629 corresponding assembly. "REAPR errors" refers to the assembly errors found by REAPR by mapping Illumina reads to the corresponding
630 assembly. For each isolate, one ONT and one PacBio-based assembly with the best completion (i.e. number of circularised contigs) were chosen
631 for comparison. DNAdiff results show the median (range) results from comparing all assemblies for an isolate across read preparation strategies
632 i.e. 4x4=16 comparisons for each isolate. "GSNPs" and "GIndels" refer to the high-confidence SNPs and indels found between ONT and PacBio
633 assemblies.

| Isolate | ALE score | PacBio unmapped reads (% total) | ONT unmapped reads (% total) | PacBio REAPR errors | ONT REAPR errors | DNAdiff GSNPs | DNAdiff GIndels |
|---|---|---|---|---|---|---|---|
| CFT073 (reference *E. coli*) | -17928 | 29246 (0.89%) | 29240 (0.89%) | 5 | 5 | 1 (0-1) | 0 (0-0) |
| MGH78578 (reference *K. pneumoniae*) | -1532602 | 41793 (1.31%) | 38371 (1.21%) | 8 | 7 | 6 (1-7) | 0 (0-1) |
| RBHSTW-00029 | 207465 | 50056 (1.85%) | 49876 (1.84%) | 3 | 3 | 0 (0-0) | 0 (0-0) |
| RBHSTW-00053 | 4727 | 50860 (1.62%) | 50861 (1.62%) | 12 | 11 | 1.5 (0-4) | 0 (0-0) |
| RBHSTW-00059 | -143627 | 37357 (1.04%) | 36251 (1.01%) | 15 | 14 | 0 (0-0) | 0 (0-0) |
| RBHSTW-00122 | 0 | 24355 (1.18%) | 24355 (1.18%) | 6 | 7 | 0 (0-0) | 0 (0-0) |
| RBHSTW-00123 | -1963188 | 56224 (1.68%) | 57074 (1.70%) | 17 | 21 | 4 (1-6) | 4.5 (2-6) |
| RBHSTW-00127 | -1145 | 34206 (0.98%) | 34206 (0.98%) | 16 | 16 | 0 (0-0) | 0 (0-0) |
| RBHSTW-00128 | 3114 | 31526 (1.06%) | 31507 (1.05%) | 6 | 8 | 2 (1-2) | 2 (1-4) |
| RBHSTW-00131 | 399368 | 25880 (0.88%) | 26271 (0.89%) | 24 | 28 | 3 (1-7) | 1 (1-3) |
| RBHSTW-00142 | -790773 | 34684 (1.23%) | 32590 (1.16%) | 12 | 12 | 3 (1-11) | 0 (0-1) |
| RBHSTW-00167 | 4083063 | 34510 (1.13%) | 76805 (2.52%) | 24 | 33 | 21 (18-47) | 1.5 (0-4) |
| RBHSTW-00189 | -158523 | 37378 (1.25%) | 37418 (1.25%) | 9 | 12 | 11.5 (7-21) | 1 (0-2) |
| RBHSTW-00277 | 18417 | 33677 (0.99%) | 33685 (0.99%) | 16 | 16 | 2 (0-2) | 0 (0-0) |
| RBHSTW-00309 | -518811 | 30704 (0.88%) | 30327 (0.87%) | 17 | 36 | 2 (0-11) | 44.5 (0-86) |
| RBHSTW-00340 | -906675 | 30802 (0.87%) | 29860 (0.84%) | 11 | 10 | 2 (0-4) | 0 (0-1) |
| RBHSTW-00350 | 21188 | 28907 (0.79%) | 28907 (0.79%) | 12 | 13 | 2 (2-4) | 5 (0-8) |
| RHB10-C07 | -23295 | 27779 (0.90%) | 27777 (0.90%) | 22 | 21 | 5 (0-17) | 0.5 (0-1) |
| RHB11-C04 | 12774 | 24879 (0.86%) | 24881 (0.86%) | 25 | 25 | 2 (0-6) | 0 (0-0) |
| RHB14-C01 | 172712 | 30478 (0.95%) | 30576 (0.95%) | 13 | 12 | 3 (0-3) | 0 (0-0) |

634 **Figure 1. Examples of genome structure uncertainty in hybrid assemblies in a) the**
635 **chromosome and b) the accessory genome.** (a) An ONT+Illumina hybrid assembly for
636 isolate RBHSTW-00029 using the "Basic" long read preparation strategy. b) A
637 PacBio+Illumina hybrid assembly for isolate MGH78578 using the "Corrected" long read
638 preparation strategy. Plots were obtained using Bandage [14], with grey boxes indicating
639 unresolved structures. Each contig is annotated with contig length and Illumina coverage;
640 connections between contigs represent overlaps between contig ends. The assembly for
641 RHBSTW-00029 in a) and that of isolate RHB14-C01 (which showed a similar pattern of
642 chromosome structure uncertainty) represented the only two datasets that could not be
643 completely assembled with any of the attempted strategies using ONT+Illumina data. They
644 were also not fully assembled by any PacBio+Illumina strategy, which similarly failed to
645 completely assemble isolates RBHSTW-00189, RBHSTW-00277, RBHSTW-340 and
646 CFT073 (Figure S4). The pattern in b) was only observed for PacBio+Illumina data, and was
647 the reason for incomplete assemblies for isolates RBHSTW-00123, RBHSTW-00131,
648 RBHSTW-00142, RBHSTW-00167 and MGH78578 (Figure S4).



(a) *C. freundii* RBHSTW-00029
ONT+Illumina, "Basic"

(b) *K. pneumoniae* MGH78578
PacBio+Illumina, "Corrected"

649

650 **Figure 2. Examples of mismatches identified between the ONT-based and the PacBio-**
651 **based assemblies for the two reference strains (*E. coli* CFT073 and *K. pneumoniae***
652 **MGH78578).** Each sub-figure is an Interactive Genome Viewer (IGV) (v2.4.3) [15] view of
653 part of the PacBio-based assembly, centered around a PacBio-ONT SNP, with all reads from
654 the same isolate mapped to it. We performed this analysis for all SNPs in isolates MGH78578
655 and CFT073, and report examples for the two most typical patterns observed. a) SNP from
656 MGH78578 with very low Illumina coverage, but normal PacBio and ONT coverage. Most
657 of the Illumina reads have a different base than the one in the PacBio-assembled reference
658 (the red T's), suggesting perhaps an error in the PacBio assembly. A similar pattern is
659 observed in 14 SNPs in CFT073 (with 12 due to error in the PacBio assembly), and 11 SNPs
660 in MGH78578 (with 10 due to error in the PacBio assembly). b) SNP from MGH78578 with
661 normal Illumina coverage; Illumina reads support both bases with similar proportions,
662 suggesting that this could be a polymorphic site within the original DNA sample. This pattern
663 was observed for 4 SNPs in CFT073 and for 13 SNPs in MGH78578.



664

## SUPPLEMENTARY FIGURES AND TABLES

**Figure S1. Read counts and read length distributions for ONT and PacBio outputs.**

**Figure S2. Summary of read-to-assembly alignments.** All assemblies considered were obtained using all reads of the given type. Reads are classified as "good" if they have at least one mapping covering 97% of the read. They are classified as a putative "chimera" if they have multiple inconsistent alignments with at least 10% of read length and 70% identity. Complete statistics from minimap2/Filtlong outputs are in Table S2.

**Figure S3. Mean percent identities and identity N50 values of ONT/PacBio reads aligned to the hybrid assemblies.** We considered the average identity for each base, and if there were multiple alignments at a base, we used the one with the best score. We aligned PacBio reads to the hybrid assembly obtained from all PacBio reads. We aligned ONT reads to the hybrid assembly obtained from all ONT reads. Identity N50 represents the percent identity for which half of the total bases are in reads with this identity value or higher. Complete statistics are in Table S3.

**Figure S4. Bandage plots for hybrid assemblies.** Each square represents one genome assembly. Shown are the ONT+Illumina (left) and PacBio+Illumina (right) assemblies for each isolate (4 columns of 5 isolates). All assembly plots are for the globally optimal long read preparation strategy for each sequencing approach i.e. "Subsampled" for ONT+Illumina and "Basic" for PacBio+Illumina (see Methods). Sequential colours for plasmids are for identical structures within isolates, but not between.

**Figure S5. Percentage of proteins with a length <90% of top UniProt hit.** Proteins in assemblies were annotated with Prokka then blasted with DIAMOND against the full UniProt database (see Methods). The proportion of proteins with a length <90% of their top UniProt hit gives a simple test for artificially shortened proteins due to indel errors in assembly. The black dashed line indicates the percentage in an existing high-quality reference genome for *E. coli* MG1655 (157 proteins out of 4240; RefSeq GCF_000005845.2). Absolute numbers were all <250; shown here is the value as a percentage of the maximum number of proteins observed in any assembly for the sample to allow comparison between different genome sizes.

**Figure S6. Comparison of discrepancy in total Prokka annotated regions across all assemblies.** The discrepancy is the number of annotated regions in the ONT+Illumina assembly minus the number of annotated regions in the PacBio+Illumina assembly. All 4x4=16 comparisons of read preparation strategies are shown.

**Table S1. Summary of sequenced isolates, DNA inputs and raw sequencing metrics.** Statistics in this table refer to raw (i.e. unfiltered) sequencing data. ONT read statistics were generated with nanostat (v0.22).

**Table S2. Classification of long reads from PacBio and ONT.** "PB" indicates PacBio. "PB2ONT" represents PacBio reads mapped to the ONT hybrid assembly, and so on. All assemblies considered were obtained using all reads of the given type. We show the number of reads falling in different categories according to how they map to the assemblies. Reads are classified as "Good" if they have at least one mapping covering 97% of the read. They are classified as a putative "chimera" if they have multiple inconsistent alignments with at least 10% of read length and 70% identity.

**Table S3. Properties of long reads from PacBio and ONT.** "PB" indicates PacBio. Reads were mapped to the assemblies using minimap2 to determine identity. We considered the average identity for each base, and if there were multiple alignments at a base, we used the one with the best score. We aligned PacBio reads to the hybrid assembly obtained from all PacBio reads. We aligned ONT reads to the hybrid assembly obtained from all ONT reads. N50 represents the length or identity for which half of the read bases are in reads of at least such length or identity.

**Table S4. Assembly runtimes in minutes.** All assemblies were run with dual 8-core Intel IvyBridge 2.6GHz, 256GB 1866MHz memory. Times include running times for Canu correction and read filtering.

**Table S5. Location and counts of proteins found uniquely in (a) ONT-based or (b) PacBio-based assembly for each sample.** Shown here is the comparison between assemblies using the globally optimal long read preparation strategy for each sequencing approach i.e. "Subsampled" for ONT+Illumina and "Basic" for PacBio+Illumina (as in Figure S4). Proteins from assemblies for each sample were clustered using Roary after annotation with Prokka. Contig order indicates size order in the relevant assembly (see Figure S4). The start of the greyed-out squares indicates the total number of contigs in the assembly.

**Table S6. Results of DNAdiff comparison between reference genomes (*E. coli* CFT073 and *K. pneumoniae* MGH78578 genomes) and hybrid assemblies with either PacBio or ONT.** Each row corresponds to a comparison, either between the reference and PacBio assembly, or between the reference and the ONT assembly, or between the two *de novo* hybrid assemblies. "Length difference" means the difference in total length of the two genomes. "Aligned bases (ref)" represents the number of bases from the first comparison genome that are aligned with the other genome in the comparison. In each comparison the ONT assembly is the one obtained using half of the long reads, while the PacBio assembly is obtained following long read correction.

**Table S7. Simulating the effect of increased level of ONT multiplexing on hybrid assembly.** Values represent numbers of contigs, either circular contigs, or any contig. Three simulations are presented, either with all reads, with half the reads, or with one third of the reads.