

1 Analysis procedures for assessing recovery of high quality,
2 complete, closed genomes from Nanopore long read
3 metagenome sequencing

4 Krithika Arumugam^{1,†}, Irina Bessarab^{2,†}, Mindia A. S. Haryono², Xianghui Liu¹,
5 Rogelio E. Zuniga–Montanez^{1,3}, Samarpita Roy¹, Guanglei Qiu^{1,4}, Daniela I.
6 Drautz–Moses¹, Ying Yu Law¹, Stefan Wuertz^{1,5}, Federico M. Lauro^{1,6}, Daniel H.
7 Huson^{7,8}, and Rohan B. H. Williams^{2,§}

8 ¹Singapore Centre for Environmental Life Sciences Engineering, Nanyang Technological
9 University, Singapore, 637551

10 ²Singapore Centre for Environmental Life Sciences Engineering, National University of Singapore,
11 Singapore, 117456

12 ³Department of Civil and Environmental Engineering, One Shields Avenue, University of
13 California, Davis, California, U.S.A. 95616

14 ⁴Current address: School of Environment and Energy, South China University of Technology,
15 Guangzhou 510006, China

16 ⁵School of Civil and Environmental Engineering, Nanyang Technological University, Singapore,
17 639798

18 ⁶Asian School of the Environment, Nanyang Technological University, Singapore, 639798

19 ⁷Center for Bioinformatics, University of Tübingen, Tübingen, Germany, 72076

20 ⁸Life Sciences Institute, National University of Singapore, Singapore, 117456

21 [†]K.A and I.B made equal contributions to this work (equal first authors)

22 [§]Corresponding author

23 March 13, 2020

24 Address for correspondence:

25 Rohan B. H. Williams, Ph.D

26 Singapore Centre for Environmental Life Sciences Engineering (SCElse)

27 Centre for Life Sciences, National University of Singapore

28 28 Medical Drive #05–02

29 Singapore 117456

30 Ph: +65–6516–4032

31 Email: lsirbhw@nus.edu.sg

Abstract

32
33 New long read sequencing technologies offer huge potential for effective
34 recovery of complete, closed genomes from complex microbial communities.
35 Using long read (MinION) obtained from an ensemble of activated sludge
36 enrichment bioreactors, we 1) describe new methods for validating long
37 read assembled genomes using their counterpart short read metagenome
38 assembled genomes; 2) assess the influence of different correction
39 procedures on genome quality and predicted gene quality and 3) contribute
40 21 new closed or complete genomes of community members, including
41 several species known to play key functional roles in wastewater bioprocesses:
42 specifically microbes known to exhibit the polyphosphate- and glycogen-
43 accumulating organism phenotypes (namely *Accumulibacter* and *Dechloromonas*,
44 and *Micropruina* and *Defluviicoccus*, respectively), and filamentous
45 bacteria (*Thiothrix*) associated with the formation and stability of
46 activated sludge flocs. Our findings further establish the feasibility of
47 long read metagenome-assembled genome recovery, and demonstrate the
48 utility of parallel sampling of moderately complex enrichment communities
49 for recovery of genomes of key functional species relevant for the study
50 of complex wastewater treatment bioprocesses.

51 The development of long read sequencing technologies, such as the Oxford Nanopore
52 Technology MinION and Pacific Biosciences SMRT are presenting new opportunities
53 for the effective recovery of complete, closed genomes [1, 2]. While these
54 new approaches have been mostly applied to single species isolates [3, 4],
55 the ability of this new methodology to recover genomes of member taxa
56 from complex microbial communities (microbiome) data is now actively
57 being explored.

58 After long read sequencing technologies first became available, several studies
59 pioneered the collection of long read data, or combined long and short
60 read data, on complex microbial communities, for example from moderately
61 to highly enriched bioreactor communities [5, 6], co-culture enrichments
62 [7], marine holobionts [8] or from full scale anaerobic digester
63 communities [9], as well as several datasets which provided benchmarking
64 data from long and short read sequencing of mock communities [10, 11,
65 12]. New long read analysis methods [13, 14] and binning algorithms
66 designed for long read metagenome data [15] have also appeared,
67 anticipating the future expansion of metagenome data generated from
68 these new instruments. More recent studies [16, 17, 18, 19, 20, 21,
69 22] have collectively established that full length (or near full length
70 genomes) can be recovered from long read sequencing of complex
71 communities, which sets the stage for further development of genome-
72 resolved long read metagenomics.

73 Here we extend our previous work [16, 22] on recovering metagenome-
assembled genomes from long read data obtained from enrichment (continuous
culture) reactors inoculated with activated sludge microbial communities.
Enrichment reactor communities [23, 24] offer a moderate level of
complexity compared to the inoculum communities [25] and so are
realistic, yet tractable, systems to use for developing approaches for
recovery and vali-

74 dation of MAG analysis using long read data. We report results and methodology of long
75 read sequencing from multiple sets of reactor communities. We have obtained short read
76 metagenome data (Illumina) from either the same DNA aliquots as used for long read
77 sequencing, or the same biomass. Specifically we 1) describe new methods for validating
78 long read assembled genomes using their counterpart short read metagenome assembled
79 genomes; 2) assess the influence of different correction procedures on genome quality and
80 predicted gene quality and 3) contribute 21 new closed or complete genomes of commu-
81 nity members, including several species known to play key functional roles in wastewater
82 bioprocesses.

83 **Methods**

84 **Overview, biomass and data availability**

85 We employed the biomass from a series of enrichment reactor microbial communities, each
86 from activated sludge sourced from wastewater treatment plants located in Singapore. We
87 sampled the following enrichment reactor communities:

- 88 (i) A lab-scale sequencing batch reactor, inoculated with activated sludge from a full
89 scale wastewater treatment plant (Public Utilities Board, Singapore), was operated
90 using acetate as the primary carbon source to enrich for polyphosphate accumulating
91 organisms (PAO). The reactor was sampled on day 267 of the operation, with both
92 long read (Nanopore MinION) and short read (Illumina Miseq 301bp PE) sequencing
93 data from the same DNA aliquot. These data have been previously published by us
94 [16] and are available via NCBI BioProject accession PRJNA509764. This data set is
95 referred to below as the *PAO1* data.
- 96 (ii) A lab-scale sequencing batch reactor, inoculated with activated sludge from a full
97 scale wastewater treatment plant (Public Utilities Board, Singapore), was operated
98 using alternative carbon sources to enrich for polyphosphate accumulating organisms
99 (PAO). The reactor was sampled on April 6, 2018, gDNA extracted and both long
100 read (Nanopore MinION) and short read (Illumina HiSeq2500 251bp PE) data ob-
101 tained from the same DNA aliquot. These data are available available via NCBI
102 BioProject accession PRJNA611629. This data set is referred to below as the *PAO2*
103 data.
- 104 (iii) Enrichment targeting putative PAO species, namely members of genera *Tetrasphaera*
105 and *Dechloromonas*. Following inoculation with activated sludge from a full scale
106 wastewater treatment plant (Public Utilities Board, Singapore), the reactor was fed
107 with synthetic wastewater containing either glutamate or glucose as the main carbon
108 source, with the feed type switched in a weekly manner, and operated at 31 °C.
109 Short read data (Illumina HiSeq2500 251bp PE) had been previously obtained from

110 sampled biomass on days 272, 279 and 286 of operation, and long read data (Nanopore
111 MinION) obtained from samples taken on days 264 and 293 of operation. These data
112 are available via NCBI BioProject accession PRJNA606905. The long read data
113 obtained from each sampling day is referred to below as the *PAO3A* and *PAO3B*
114 data, respectively.

115 (iv) Enrichment targeting PAO species capable of performing denitrification. A lab-
116 scale sequencing batch reactor was inoculated with activated sludge from a full-
117 scale wastewater treatment plant (PUB, Singapore). The reactors were operated at
118 35 °C using acetate as the primary carbon source fed under anaerobic conditions,
119 but without the addition of allyl-thiourea (ATU) in order to suppress the growth of
120 ammonia oxidizing bacteria, with the aim of targeting polyphosphate-accumulating
121 organisms that could also reduce nitrogen oxides (nitrite and/or nitrate). For this
122 study we obtained long read data (Nanopore MinION) and short read data (Illumina
123 HiSeq2500 251bp PE) from the same DNA aliquot extracted from biomass sampled
124 on day 292 of operation. These data are available via NCBI BioProject accession
125 PRJNA607349. These data are referred to below as the *PAO4* data.

126 Using these data, our main objective was to obtain complete bacterial chromosomes, via
127 assembly of long read data, and use draft genomes obtained from metagenome assembly
128 of corresponding short read data for the purposes of evaluation. This approach takes
129 advantage of the fact of our having obtained data from both sequencing modalities and
130 takes advantage of current understanding of short read metagenome assembly binning and
131 quality assessment procedures [26, 27]. We highlight that it may be possible to recover
132 further genomes from these data by the use of binning procedures adapted to long read
133 data, however here our focus is on analysing contigs directly obtained from the assembly
134 that plausibly represent complete bacterial chromosomes.

135 DNA extraction

136 Genomic DNA in the case of the samples from PAO1, PAO3A, PAO3B and PAO4 was
137 extracted from the sampled biomass as described previously by us [16], briefly, we used the
138 FastDNATMSPIN Kit for Soil (MP Biomedicals), using 2× bead beating with a FastPrep
139 homogenizer (MP Biomedicals). Extracted gDNA from the PAO1, PAO3A and PAO3B
140 samples was then size-selected on a BluePippin DNA size selection device (SageScience)
141 using a BLF-7510 cassette with high pass filtering with a 8 kbp cut-off. The gDNA
142 from the PAO4 sample was size-selected using Circulomics Short Read Eliminator XS kit
143 (Circulomics Inc). Size-selected DNA was then taken for Nanopore library construction
144 (see below).

145 From the biomass from PAO2 sampling, high molecular weight (HMW) DNA was
146 extracted using a modified xanthogenate-SDS protocol [28]. Briefly, 2 mL of biomass
147 from lab-scale sequencing batch reactor was harvested by centrifugation, the pellet was

148 resuspended in 0.6 mL of DNA/RNA shield (Zymo Research) and added to 5.4 mL of pre-
149 heated (65 °C) XSP buffer (1:1 volumes of XS buffer and phenol). The tubes were incubated
150 at 65 °C for 15 min, vortexed for 10–15 sec, placed on ice for 15 min and centrifuged at
151 14000 rpm for 5 min. The aqueous phase was transferred to a fresh tube and extracted with
152 equal volume of phenol:chloroform:isoamyl alcohol (25:24:1) followed by extraction with
153 chloroform:isoamyl alcohol (24:1). The aqueous phase after chloroform:isoamyl alcohol
154 (24:1) extraction was ethanol precipitated and resuspended in TE buffer. The extracted
155 DNA was further treated with RNase A (Promega) then extraction with phenol, followed
156 by phenol:chloroform:isoamyl alcohol (25:24:1), and ethanol precipitation. Purified DNA
157 was taken to library construction for Nanopore sequencing.

158 **Short read sequencing**

159 Genomic DNA Library preparation was performed using a modified version of the Illumina
160 TruSeq DNA Sample Preparation protocol. We then performed a MiSeq sequencing run
161 with a read length of 301 bp (paired-end) or a HiSeq2500 sequencing run with a read
162 length of 251 bp (paired-end) as specified above.

163 **Long read sequencing**

164 Nanopore sequencing was performed on a MinION Mk1B instrument (Oxford Nanopore
165 Technologies) using a SpotON FLO MIN106 flowcells and R9.4 chemistry. Data acqui-
166 sition was performed using MinKNOW software, without live basecalling, running on a
167 HP ProDesk 600G2 computer (64-bit, 16 GB RAM, 2 Tb SSD HD; Windows 10). The
168 runs were continued until active pores in flowcells were depleted. For PAO1, PAO3A and
169 PAO3B extractions, the sequencing library was constructed from approximately 4–4.5 μg
170 of size-selected genomic DNA using SQK-LSK108 Ligation Sequencing Kit and approxi-
171 mately 900 ng of the library was loaded onto each flowcells. For PAO2 data set, sequencing
172 libraries were constructed from HMW DNA using two different sequencing kits from ONT.
173 The first kit was the Rapid Sequencing kit SQK-RAD004, for which the library was con-
174 structed from 400 ng of HMW DNA and the entire library loaded onto the flow cell. The
175 second kit was the Ligation Sequencing kit SQK-LSK 108, for which 1.0 μg of genomic
176 DNA was used for library construction, and 400 ng was loaded onto the flow cell. For
177 PAO4 data set, the sequencing library was constructed from 1.2–1.3 μg of size selected
178 DNA using SQK-LSK109 Ligation Sequencing Kit (Oxford Nanopore Technologies). The
179 library was diluted to allow 250 ng of the library to be loaded on the flowcell.

180 **Analysis of long read sequence data**

181 Basecalling was performed with **guppy** (CPU version 3.2.1, 3.2.2 or 3.3.0 for Linux 64-bit
182 machines; see **Table S1**). Adaptor trimming was performed using **Porechop** (version 0.2.2)

183 [29] with default settings except `-v 3 -t 20`. We assembled long read data using `Canu` (ver-
184 sion 1.8 or 1.9, default settings except `corOutCoverage=10000`, `corMhapSensitivity=high`,
185 `corMinCoverage=0`, `redMemory=32`, `oeaMemory=32` and `batMemory=200 useGrid=false`)
186 [30], `Unicycler` (version 0.4.7 or version 0.4.8 with default settings except `-t 20 --keep`
187 `3`) [31] and `Flye` (version 2.4 with default settings except `-t 20 --plasmids --debug`
188 `--meta`) [32]. Contigs generated from long read data are hereafter denoted as *long read*
189 *assembled contigs* (LRAC). The number of reads used in each assembly was estimated
190 by mapping long read to LRAC sequence with `minimap2` (version 2.17) [33] and using
191 `samtools-1.6` to calculate the number of aligned reads [34]. We used `DIAMOND` (version
192 0.9.24) [35] to perform alignment of LRAC sequences (with default settings except `-f`
193 `100 -p 40 -v --log --long-reads -c1 -b12`) against the NCBI-NR database (Febru-
194 ary, 2019) [36]. From the MEGAN Community Edition suite (version 6.17.0) [37] we
195 used `daa-meganizer` (run with default settings except `--longReads`, `--lcaAlgorithm`
196 `longReads`, `--lcaCoveragePercent 51`, `--readAssignmentMode alignedBases` and the
197 following settings for mapping files: `--acc2taxa prot_acc2tax-Nov2018X1.abin`, `--acc2eggnog`
198 `acc2eggnog-0ct2016X.abin`, `--acc2interpro2go acc2interpro-June2018X.abin`, `--acc2vfdb`
199 `acc2vfdb-Feb2019.map`) to format the `.daa` output file for use in the MEGAN GUI (ver-
200 sion 6.17.0). Within MEGAN, LRAC sequences were exported with the ‘Export Frame-
201 Shift Corrected Reads’ option to obtain frameshift corrected sequence. LRAC sequence
202 that was at least 1Mb in length were categorised as potential whole chromosome sequence
203 and from thereon described as *LR-chr* sequence. We processed LR-chr sequences with
204 `CheckM` (version 1.0.11) [38] and `Prokka` (version 1.13) [39] to assess genome quality. LR-
205 chr sequences that demonstrated `CheckM-SCG` completeness 90% and contamination < 5%
206 were classified as putative genomes. The entire set of putative genomes were derepli-
207 cated using the `dRep` (version 2.2.3) workflow [40] with the following settings: `-p 44`
208 `-comp 90 -con 5 -str 100 --genomeInfo`. We performed taxonomic annotation on re-
209 covered genome sequence using `GTDB-Tk` (version v0.3.2, running default parameters ex-
210 cept `--cpus 40 -x fasta`) [41]. Coverage profiles were generated from both long read
211 and short read data, by mapping each of these to LR-chr sequences using `minimap2` (ver-
212 sion 2.17) using the following flags `-ax map-ont` for long read data, and `-ax sr -a -t`
213 `20` for short read data. Sorted `.bam` files were subsequently processed using `bedtools`
214 `genomeCoverageBed` (version 2.26.0) with the following flags `-d`. We extracted 16S-
215 SSU rRNA genes as identified with `Prokka` and annotated them against SILVA database
216 (SUrRef_NR99_132_SILVA_13.12.17_opt.arb) [42] using `sina-1.6.0-linux` [43] running de-
217 fault settings except `-t -v --log-file --meta-fmt csv` and with `--lca-fields` set for
218 all five databases, namely `tax_slv`, `tax_embl`, `tax_ltp`, `tax_gg` and `tax_rdp`.

219 Analysis of short read sequence data

220 The raw FASTQ files were processed using `cutadapt` (version 1.14) [44] with the fol-
221 lowing arguments: `--overlap 10 -m 30 -q 20,20 --quality-base 33`. We performed

222 metagenome assemblies from short read data using SPAdes (version 3.12.0-Linux or 3.14.0-
223 Linux, with default settings except `-k 21,33,55,77,99,127 --meta`) [45] either as single
224 sample assemblies, in the case of short read data from PAO1, PAO2 and PAO4, or as
225 co-assembly of all short read in the case of the PAO3A and PAO3B samples. The con-
226 tigs generated from short read data are hereafter denoted as *short read assembled con-*
227 *tigs* (SRAC). We identified putative member genomes using MetaBAT2 [46], after filtering
228 for contigs at least 2000 bp in length. We identified 16S genes within contigs using the
229 `--search16` module of USEARCH (version 10.0.240, 64 bit) [47], and annotated them using
230 the SILVA::SINA webserver (using default parameters) [42, 43]. For each identified bin we
231 performed genome quality estimation using CheckM (version 1.0.11). We performed tax-
232 onomic annotation on recovered member genomes using GTDB-Tk (version 0.3.2, running
233 default parameters except `--cpus 40 -x fasta`).

234 Comparative analysis of long and short read assemblies

235 We used BLASTN (version 2.7.1+) [48] to examine the degree of sequence alignment between
236 LRAC and SRAC sequences. We treated the LRAC as the subject sequences and the SRAC
237 as the query sequences, using default BLASTN parameters (except `-outfmt 6`). From the
238 BLASTN tabular output, we retained the highest bit-score from each unique combination of
239 query and subject pair. In order to identify the short reads bin(s) that are cognate to a given
240 LR-chr sequence, we then computed the *concordance statistic* (κ), as previously described
241 by us in [16], for all combinations of short read contigs (categorised by bin membership) and
242 LRAC sequence, that were present in the BLASTN output. We then compute the following
243 component statistics:

- 244 (1) \widehat{pid} : The mean of the percent identity (PID), calculated across alignments, and
245 quantified as a proportion. \widehat{pid} is defined on the interval [0,1]
- 246 (2) $\widehat{al2ql}$: The mean of the quotient of the alignment length to the query length, cal-
247 culated across alignments, and quantified as a proportion. $\widehat{al2ql}$ is always ≥ 0 and
248 while values > 1 can be observed, in practice the maximum observed value be ap-
249 proximately 1.
- 250 (3) p_{srac} : the quotient of the number of short read contigs in the bin that produce
251 alignments and the total number of short read contigs in the bin. p_{srac} is defined on
252 interval [0,1]
- 253 (4) p_{aln} : the proportion of the long read contig that is covered by an alignment. p_{aln} is
254 defined on the interval [0,1].

255 Collectively these statistics contain information on how well a set of short read contigs will
256 tile a LR-chr sequence, namely, completeness of coverage (captured by p_{srac} and p_{aln}), as
257 well as quality of the alignments (captured by \widehat{pid} and $\widehat{al2ql}$). We can hypothesise that if

258 the majority of the contigs in a short read MAG completely covered a LR-chr sequence
259 with high quality alignments, we would predict all four of these statistics would hold values
260 be close to unity. A simple extension of this prediction is to calculate the mean of the four
261 statistics, which we denote as the concordance statistic, $\kappa = (\widehat{pid} + \widehat{al2ql} + p_{srac} + p_{aln})/4$,
262 which provides a single number to screen large numbers of pairwise combinations of short
263 read and long read derived MAG in an efficient way. The concordance statistic (κ) was
264 computed using all alignments, as well as after filtering for near-full length alignments
265 (defined as $al2ql \geq 0.95$). We provide an R package `srac2lrac` to compute κ (along
266 with all component statistics) following calculation of BLASTN-like alignment statistics and
267 definition of short read bins.

268 **Analysis of effects of frame-shift correction on coding sequence**

269 The frameshift correction procedures employed from the MEGAN-LR package [13, 16] have
270 been crucial in permitting the use of genome quality and annotation workflows, namely
271 `CheckM` and `Prokka`, and we further evaluated the extent to which these procedures im-
272 proved accuracy of the coding gene sequence. To do so we analysed the distribution of
273 ratio of predicted gene length to the length of the nearest orthologous gene, as suggested
274 by Watson and colleagues [49], before and after the application of frameshift correction
275 procedures, as well as against two other sequence correction algorithms, namely `Medaka`
276 (version 0.11.5) [50] and `Racon` (version 1.4.3) [51]. In the first instance, we employ a single
277 round of correction and in the case of `Racon` do not use short read data for correction (in
278 order to maintain independence of each data type, as in the case of the concordance statis-
279 tic calculations). MEGAN-LR frameshift correction procedure uses the results of alignments
280 made against RefSeq NR, we did not compare against this same database to avoid positively
281 biasing the performance, rather we used predicted genes from the cognate short read assem-
282 blies as a database of genes to use as subject sequences. Specifically, we took the protein
283 coding sequence of each ORF in each of four versions of the genome generated above, and
284 performed homology search of each sequence against the short read assembly ORF database
285 using `DIAMOND` (version 0.9.24, running in `blastp` mode with default parameters except `-f`
286 `6 qseqid qlen slen sseqid sallseqid -p 40 -v --log --max-target-seqs 1`). We
287 then calculated the quotient of the length of the query sequence to the length of subject
288 sequence holding the maximum bit-score (best hit). We ran `CheckM` on each of the four
289 versions of each putative genome, as described above. We also examined the common prac-
290 tice of applying multiple rounds of correction, with within and across, different correction
291 software, by performing both the above analyses on genomes corrected with four sequential
292 applications of `Racon` followed by one application of `Medaka` (denoted as ‘multiple’ from
293 hereon).

294 Procedures for refining draft genomes

295 In LR–chr sequence we screened regions of potential misassembly by identifying genomic
296 intervals of at least 10bp in length, where long read coverage was either abnormally high
297 or abnormally low, defined as >1.5 of the median coverage and <0.5 of the median cover-
298 age, respectively. We then examined alignments of both long and short read data to the
299 genomes using the Integrated Genome Viewer (IGV version 2.4.14) [52] to identify low cover-
300 age regions that showed evidence of misconnection between reads, or weakly supported
301 connection, or in the case of high coverage regions, to disambiguate types of read con-
302 nections likely to arise from non-cognate sources. We generated VCF files for short read
303 alignments using `BCFtools` (version 1.9 run with flags `-mv`) [53] to identify likely single
304 nucleotide variants and presence of insertion/deletion variants, and subsequently used the
305 aligned short read contig sequences to remove false nucleotide calls. We then align the
306 entire genome against itself using BLASTN to check the integrity of the corrected genome
307 sequence. For completeness, we have provided raw LR–chr sequence, frame-shift–
308 corrected sequence and, for the subset of genomes subjected to further refinement, the
309 fully completed versions.

310 Data availability

311 Raw sequence data is available at NCBI Short Read Archive (SRA) via BioProject accession
312 identifiers listed above. The R code for performing the concordance statistic analysis are
313 available at <https://github.com/rbhwilliams/srac2lrac> including test data and scripts
314 taken from the PAO2 data.

315 A Zenodo submission (<https://doi.org/10.5281/zenodo.3695987>) contains key sec-
316 ondary data, including: 1) LRAC sequence from each dataset; 2) whole genome sequence
317 from the 21 genomes listed in **Table 1** for each of the five correction procedures (FASTA se-
318 quence, `Prokka` and `CheckM` results); 2) short read assembled sequence and binning results;
319 3) concordance statistic data and results; 4) short and long read per–base coverage data
320 for the 21 genomes and 5) two manually corrected genomes of *Candidatus Accumulibacter*
321 (also being submitted to NCBI) along with detailed notes explaining the procedures that
322 were applied.

323 Results

324 Long read sequencing depth improved from the beginning of the study period, reflecting
325 rapid improvement in experimental protocols and flow cell technology (**Table S1**), with
326 the total amount of sequence generated ranging from around 1Gbp/run to just under
327 12Gbp/run (**Table S1**). These data were assembled using each of the three workflows
328 as described above. The Canu assembly workflow generated a greater number of LR–
329 chr ($n=90$) on these data than did either Unicycler ($n=44$) or Flye ($n=60$) (**Table S2**).

330 As Canu generated a substantially larger number of LR-chr sequences, we subsequently
331 focused attention on the results obtained with this workflow (see **Table S3** for comparative
332 summary of LR-chr sequences from each workflow).

333 We next applied the truncated MIMAG criteria for estimating high quality MAG status
334 (SCG-estimated completeness > 90% and contamination < 5%) and observed a total of
335 23 LR-chr generated from Canu that could be considered plausible candidates for being
336 whole chromosomal sequence, from here on referred to as *putative genomes* for convenience.
337 A further 13 LR-chr sequences from Canu were classifiable as medium quality (SCG-
338 estimated completeness \geq 50% and contamination < 10%, including one that was LR-chr
339 was circular). We de-replicated the entire set of 23 putative genomes using the dRep
340 workflow with a relatedness threshold of ANImf>99 (**Figure S1**), obtaining a reduced set
341 of 21 putative genomes (**Table 1**). The two redundant genomes were obtained from the
342 PAO3A and PAO3B datasets, consistent with the fact that they are the same community
343 sampled at different times.

344 We then studied each of these 21 dereplicated putative genomes in more detail to
345 establish whether they were, or were not, likely to represent whole chromosomes. Using
346 annotations from the Prokka workflow, all 21 putative genomes met the complete MIMAG
347 criteria for being classified as high quality metagenome assembled genomes, including a
348 minimum number of tRNA encoding genes, and the presence of each of the genes encoding
349 5S, 16S and 23S SSU-rRNA genes detected in each genome. Estimated SCG completeness
350 was on average 95.87% (range 93.47–99.04%) and mean contamination was 0.37% (range
351 0.00–1.09%). Nine of the 21 sequences were classified as circular by Canu (**Table 1**).
352 Coverage profiles generated using both long and short read data within a given community
353 showed uniform coverage, with no substantive gaps observed (see Panel C of **Figure 1** and
354 **Supplementary Figures 2–30**). The proportion of long reads utilised to produce the
355 putative genomes varied with dataset (**Table S4**) but from conservative estimation (based
356 on subsetting alignments of all long reads against all LRAC sequence), on average 32.6%
357 of reads across all 5 data sets (range: 23.3–55.4). At the individual genome level as few as
358 1% of reads in a dataset could generate a complete genome (PAO1-tig00000003; see **Table**
359 **S4**).

360 To gain further insight into the quality and completeness of detected genomes, we used
361 the *concordance statistic* (κ), previously developed by us [16] to identify metagenome-
362 assembled genomes obtained from short read sequence data that were cognate to a long
363 read assembled genome (summary data of each short read assembly is provided in **Table**
364 **S5**). The κ -statistic is computed for all combinations of short read MAG and LR-chr
365 sequences. A observed value of κ close to unity will imply that the LR-chr sequence is tiled
366 by the contigs from the short read MAG, and the latter can be considered a likely candidate
367 for being the cognate genome. For 20 of the 21 genomes in **Table 1** the maximum observed
368 κ values were high (mean: 0.95 range: 0.83–1.00) (**Table S6**). If we only considered near
369 full length alignments ($al2ql > 0.95$), this reduced by around 0.5 units (mean 0.89, median
370 0.91, range: 0.80–0.97). In **Figure 1** we provide a comprehensive visualisation of the

371 concordance statistic analysis for the case of the PAO3A–tig00018026 genome against its
372 cognate short read MAG (bin 114). Related plots for all 21 genomes are available in the
373 **Supplementary Figures 2–30**. We observed a genome recovered from PAO4 (PAO4–
374 tig00000079), annotated at species level to *Exiguobacterium profundum*, which held a κ
375 value of 0.3 and from which there appeared to be no corresponding complete short read
376 MAG (**Supplementary Figure 29**).

377 On average, for a given LR–chr sequence, κ –statistics were generated from around two
378 thirds of available short read MAGs, but in most cases the magnitude of the κ –statistic
379 itself was low. Of the four component statistics, \widehat{pid} and p_{srac} showed consistently higher
380 values in the bulk of associations than either $\widehat{al2ql}$ or p_{aln} , with the latter two measures
381 provided greater visual discrimination between the short read MAG holding the maximum
382 κ value and the bulk distribution of (lower) κ scores. As expected, cognate short read
383 MAGs were generally drawn from among the most abundant members of a given reactor
384 community. Contigs from short read bins with related taxonomy usually scored highly on
385 one or more component scores (data not shown), but in combination, only one short read
386 MAG generated a high value κ score with component statistics that supported it being the
387 cognate. In several cases, we observe two short MAG that tile two adjacent fragments of a
388 single LR–chr sequence, which we determined to be due to underlying genome being split
389 by MetaBat2 into two or more component sub–MAG (bin–splitting; see **Supplementary**
390 **Figures 2–3, 5–6, 10–11, 20–22**).

391 **Identification of probable mis-assemblies among LR–chr sequences**

392 Among the complete set of LR–chr we identified several examples of LR–chr that are clearly
393 mis–assemblies. In the PAO3A data, we observed one contig (tig00000001; assembled by
394 Canu) that appeared to be comprised of two separate complete genomes joined together
395 (see **Supplementary Figures 31–34** for further dissection). In this case, the proximal
396 two thirds of the LR–chr arises from one genome, while distal third from another, as
397 evidenced by different GC proportions and divergent short bin associations, respectively.
398 In the case of the PAO4 data we observed several LR–chr that were classified by CheckM to
399 have completeness over 90% but which demonstrated substantial degrees of contamination
400 (namely tig00017984, tig00017990 and tig00017987 from Canu), most likely as the results
401 of reads from closely related strains being combined.

402 **Effect of sequence error correction on coding sequence and genome quality**

403 Although the recovered genomes are consistent with being *bone fide* whole bacterial chromo-
404 somal sequence, the high error rate present in current nanopore–based sequencing implies
405 these constructs may not meet current expectations of reference genome quality. Examining
406 the length ratio histograms of the predicted genes from long read assembled genomes,
407 against their best hit counterparts from the cognate short read assemblies, we observed

408 that application of any of the three sequence correction procedures provided some degree
409 of improvement compared to the case of raw sequence, with an increased frequency of the
410 length ratio being located around a value of unity (**Figure 2** and **Supplementary Figure**
411 **35**). The performance of **Racon** was highly variable but always less effective than either
412 **MEGAN-LR** or **Medaka**. **MEGAN-LR** generally provided the best performance, followed by the
413 multiple procedure approach, than **Medaka**.

414 We further examined the influence of sequence correction on genome quality statistics,
415 as estimated by **CheckM** (**Table 2**). Of the 21 frame-shift corrected genomes classifiable as
416 high quality (**Table 1**), 3, 13, 7 and 16 of these were also classifiable as high quality when
417 examined in their uncorrected, **Medaka**-corrected, **Racon**-corrected and multiple procedure
418 corrected forms (**Table 2**), with the mean completeness being 76.5% (range: 41.9–93.0),
419 91.7% (range: 78.0–98.0), 86.8% (range: 66.0–97.0) and 92.0% (range: 77.2–98.0), respec-
420 tively, compared to a mean of 95.3% (range: 92.6–99.0) in the case of **MEGAN-LR**. Contam-
421 ination was never observed to be greater than 3% in any version of the 21 genomes.

422 Taxonomic analysis of recovered genomes

423 We inferred taxonomy of the recovered genomes using **GTDB-Tk**, as provided in **Table**
424 **S7** and summarised below (additionally we provide a complementary analysis of recovered
425 16S-SSU rRNA gene sequence annotated against the **SILVA** database in **Table S8**, and
426 **GTDB-Tk** annotations for all short reads bins in **Table S9**). Of the 21 long read genomes,
427 5 had sufficiently high degree of similarity to be classified to species level and 10 to genus
428 level, 5 to family level and 1 to class level.

429 We recovered genomes of four taxa that hold known relevance to wastewater biopro-
430 cess, namely two genomes from the PAO species *Accumulibacter*: the PAO1-tig00000001
431 genome was closely related to *Candidatus Accumulibacter* sp. SK-02 [54], and found in
432 our previous analysis of the PAO1 data [16], and the other (PAO2-tig00000001) related to
433 *Ca. Accumulibacter* sp. BA-94 [54]; **Figure 3**) and a short read MAG previously recov-
434 ered by us and denoted as *Candidatus Accumulibacter* clade IIF Strain SCELSE-1 [55].
435 We recovered two genomes for *Dechloromonas*, generally considered as exhibiting the PAO
436 phenotype [56], and one of genus *Micropruina*, previously shown to exhibit the glycogen
437 accumulating organism (GAO) phenotype [57, 58]. The PAO1-tig00026549 sequence, an-
438 notated to the novel **GTDB**-derived family 2-12-FULL-67-15 and harbouring a 16S gene
439 annotated to *Deftuviococcus* (**Table S8**), represents a novel member of the latter genus,
440 whose members exhibit the GAO phenotype [59]. We also recovered a genome from a
441 member of genus *Thiothrix*, a filamentous bacterium associated with the maintenance of
442 floccular structure in activated sludge biomass [60, 61].

443 A set of four genomes recovered here have been previously identified in temperate cli-
444 mate activated sludge, namely 3 members of the **CFB** group recovered from the PAO1 data,
445 OLB8 [62], OLB11 and OLB12 [62], as well as a genome classified to the *Rhodobacteraceae*
446 genus *UBA1943* [63]. *Thiobacillus* has been previously identified in activated sludge from

447 industrial wastewater treatment plants [64, 65]. In the PAO4 community, we recovered
448 a genome close to that of *Exiguobacterium profundum*, originally discovered in deep-sea
449 hypothermal vents [66]. Members of this genus, namely *Exiguobacterium alkaliphilum* and
450 *Exiguobacterium* sp. YS1, have been studied in relation to treatment of high alkaline
451 brewery wastewater and solubilisation of waster activated sludge, respectively [67, 68].
452 The genome of a member of family *Parachlamydiaceae*, an environmental *Chlamydia* [69]
453 that was previously recovered by us [16] in the PAO1 data, and is probably a symbiont
454 species of protists that are known to inhabit activated sludge [70]. The genome of *Bre-*
455 *vundimonas* was closely related to a short read MAG previously obtained from an activated
456 sludge metagenome in Hong Kong [71], and members of this genus have been observed pre-
457 viously in activated sludge systems [72], where they have been associated with quinoline
458 degradation from coking wastewater [73].

459 A genome from a member of genus *Pseudoxanthomonas* was also recovered. The re-
460 maining genomes had no close references, and likely represent novel members of the micro-
461 bial groups, namely family *Nocardioideaceae* (tig00157979 from PAO3B), from class *Anaero-*
462 *lineae* (tig00018026 from PAO3A), family *Burkholderiaceae* (tig00000024 from PAO3B),
463 and a genome from the novel UBA6002 family (tig00000117 from PAO1).

464 **Further refinement of genomes of *Candidatus Accumulibacter***

465 We applied manual refinement procedures, as described in **Materials and Methods**, to
466 the two recovered genomes of *Candidatus Accumulibacter*, namely PAO1–tig00000001 and
467 PAO2–tig00000001, in order to obtain submission quality finished genomes. Detailed notes
468 on refinement and manually curation procedures are provided in the Zenodo submission.

469 **Discussion**

470 In this paper we explore how long read metagenome data, generated by a Nanopore MinION
471 sequencer, can enhance the recovery of member genomes of microbial communities. Build-
472 ing on our previous analyses [13, 16, 22], we obtain further data from activated sludge
473 enrichment bioreactor communities, and obtain 21 non-redundant complete genomes, of
474 which nine are closed (circular) and six are from species with key functional relevance to
475 wastewater bioprocesses. Additionally, we present further details of methodology for as-
476 sessing whether genomes obtained from short read assemblies recapitulate those obtained
477 from assembled long read data (the concordance statistic, briefly introduced by us in [16]),
478 and examine aspects of genome quality not previously covered, including the quality of gene
479 level coding sequence and the sequence rising from the mis-assembly and related artefacts.
480 These new results highlight that by using long read sequencing in microbial communities
481 of moderate complexity, it is clearly feasible to capture sequence constructions that are
482 close to the requirements of high quality, closed genomes, for the most abundance commu-
483 nity members, without the use of contig binning procedures. However careful evaluation

484 of such genomes still appears mandatory to assess quality and the presence of artefactual
485 constructs.

486 The wide spread use of metagenome-assembled genomes (MAG) methodology on short
487 read metagenome data has provided a tremendous number of new draft genomes from di-
488 verse microbiomes and microbial communities (for example [20, 54] among others). How-
489 ever substantial limitations of these approaches have become evident, including problems
490 related to the use of multi-sample co-assemblies [20, 74, 75], the challenges of resolving
491 genomes to strain level [76], difficulties related to extracting MAGs from communities of
492 high ecological complexity [77, 78], and the limitations of automated binning procedures,
493 requiring careful evaluation of recovered genomes [79]. In response to these challenges,
494 recent efforts have combined short read with emerging complementary techniques such as
495 HiC metagenomics [80, 81, 82], synthetic long reads [83, 84], or long read sequencing, and
496 collectively these results suggest substantial improvement can be made in the quality and
497 completeness of metagenome assembled genomes using multiple types of sequence data. In
498 the present study, we make use of DNA extractions that are co-assayed with both long
499 and short sequencing, or DNA extractions from sampling events close enough together in
500 time, that we can discount the influence of the ecogenomic differences as major influence
501 on any observed differences between the types of sequence data.

502 Our analysis proceeds on the basis that neither short read nor long read data can
503 be assumed to provide an accurate reference genome, and so we seek to understand and
504 characterize the degree of agreement between assembled sequence generated from each
505 data source. Although error prone MinION sequence can be corrected using higher quality
506 short read sequences [41, 42], we have deliberately kept the two sources of data separate
507 so as to not introduce any positive bias in the calculation of the concordance statistics.
508 The concordance statistic was developed to provide a straightforward screening procedure
509 for identifying short read MAG that are cognate to assembled genomes from long read
510 data, by capturing information from alignment statistics. The concordance statistic also
511 may have broader utility, for example we can observe several instances of 'split' bins from
512 the short read assembly that are cognate to a given long read assembled genome, and
513 cases where assembled long read sequence is demonstrably artefactual. We highlight that
514 the concordance statistics capture more information than are contained in dot-plots (which
515 require the imposition of arbitrary decision thresholds on alignment statistics), albeit at the
516 cost of increased complexity. We provide R code for computing the concordance statistics
517 from alignment results, and example workflows for visualisation.

518 While these are clearly vast improvements on the working models of genomes available
519 from short read MAG analysis, several problems are still present that require attention
520 and/or explicit correction. Firstly, the high error rate implicit in MinION sequence (not
521 less than 5% sequencing error [85]) requires correction procedures to be applied either pre-
522 or post-assembly. In the present case, we are relying on the frameshift correction algorithm
523 implemented in MEGAN-LR [13, 16], which appears to perform slightly better than the
524 next best correction procedure ('multiple'). As previously discussed [16], this correction

525 procedure permit the application of existing genome quality workflows (**CheckM** in the
526 present case), and the resulting sequence can be considered to be at least a high quality
527 assembly under currently accepted criteria (MIMAG as defined in [27]). However further
528 analysis of the corrected gene content suggests there remains a substantial proportion of
529 genes that remain inadequately corrected when compared against genes predicted from the
530 cognate short read assemblies. Because the **MEGAN-LR** correction is dependent on aligned
531 sequence from database comparisons, a combination of false positives alignments and a lack
532 of closely related reference genomes could result in inappropriate or inadequate correction
533 of the query sequences in our analysis, and additionally mis-assembly of genes in the
534 short read assembly (subject sequences) could also be a factor in patterning these findings.
535 A second factor relates to the inclusion of artefactual sequence (mis-assembly), which we
536 identify and remove using examination of read alignment and coverage profiles, in line with
537 recent calls for the continuing need for careful evaluation of the output of automated genome
538 recovery procedures [79]. Collectively these results indicate that long read sequencing
539 technology that harbours high error rates should be considered complementary to short
540 read sequencing for the foreseeable future, with self-evident implications for experimental
541 design choices.

542 We have deliberately focused on long read assembled contigs that form single contiguous
543 sequences that are consistent with being whole bacterial chromosomes, which plays to the
544 full strengths of long read sequencing. The remaining, much larger set of contigs, that
545 do not meet our criteria for being considered putative genomes, will be in part comprised
546 of genome fragments that could be recovered into draft genomes using binning methods.
547 Although the amount of long read metagenome data collected from microbial communities
548 of high to very high complexity is only just emerging [21], recent work on combining short
549 and long read data from human fecal microbiomes [19] suggests that binning procedures
550 will have to be developed, or adapted from short read methods, for the full potential of
551 these new hybrid data to be realised.

552 In the present study, we are able to draw strength from the fact that the communities
553 under study are of moderate complexity, and, in ecological terms are of low evenness,
554 compared to the source inoculum, namely activated sludge residing in full scale wastewater
555 treatment plants [25]. This suggests that one way to approach a systematic genome-
556 resolved dissection of such complex communities would be to simply sample a diverse array
557 of such enrichment communities, rather than rely on more deeper, expensive sequencing
558 of a limited number of highly complex source communities. While such an approach may
559 miss some relevant species (due to the biases of enrichment protocols), it would permit the
560 recovery of many near-finished genomes from key species of direct functional relevance to
561 wastewater bioprocess engineering, as obtained here.

562 List of Supplementary Materials

563 List of Supplementary Tables (.xlsx format)

- 564 • Supplementary Table 1: Summary of long read data.
- 565 • Supplementary Table 2: Summary statistics for long read assemblies from three as-
566 ssembly workflows.
- 567 • Supplementary Table 3: Comparison of number and CheckM-derived genome quality
568 statistics of LR-chr sequences generated by three assembly workflows.
- 569 • Supplementary Table 4: Estimation of long read read count used for assembly of
570 recovered genomes.
- 571 • Supplementary Table 5: Summary statistics for short read assemblies.
- 572 • Supplementary Table 6: Summary of concordance analysis for recovered genomes.
- 573 • Supplementary Table 7: Taxonomic annotation of genomes using GTDB-Tk.
- 574 • Supplementary Table 8: Taxonomic annotation of genomes from 16S sequence.
- 575 • Supplementary Table 9: Taxonomic annotation of short read MAG using GTDB-Tk.

576 Supplementary Figures

- 577 • Supplementary Figure 1: Tree generated from MASH distances between 23 LR-chr
578 classifiable as putative genomes and used to undertake genome dereplication.
- 579 • Supplementary Figures 2–30: Summary of concordance statistic analysis for all 21
580 genomes listed in Table 1.
- 581 • Supplementary Figures 31–34: Summary of concordance statistic analysis for an
582 artefactual LR-chr sequence (PAO3A-tig00000001).
- 583 • Supplementary Figures 35: Figure 2 presented with a logarithmic scale on the vertical
584 axis.

585 Author contributions

586 The study was designed by R.B.H.W and I.B. R.E.Z.M, S.R, G.L.Q, Y.Y.L and S.W setup
587 and operated enrichment reactors, and obtained samples with I.B. I.B and F.L designed
588 long read sequencing experiments and I.B performed DNA extractions and performed long
589 read sequencing. D.I.D-M obtained short read sequencing data. K.A, M.A.S.H, D.H.H.,

590 X.H.L and R.B.H.W designed analyses, performed data analysis and/or wrote analysis
591 code. All authors contributed to data interpretation. R.B.H.W wrote the manuscript with
592 specific contributions from all other authors. K.A and I.B made equal contributions to this
593 work.

594 Acknowledgements

595 This research was supported by the Singapore National Research Foundation and Ministry
596 of Education under the Research Centre of Excellence Programme and by program grant
597 1301-IRIS-59 from the National Research Foundation (NRF). The computational work was
598 performed in part on resources of the National Supercomputing Centre (NSCC) supported
599 by Project 11000984. We thank Gavin Huttley (Australian National University) for critical
600 feedback on sequence analysis, Uma Shankari d/o Chanda Segaran for performing nucleic
601 acid co-extraction from the PAO3 samples and constructive, critical reviews from an earlier
602 submission which has vastly improved this paper.

603 References

- 604 [1] Loman, N.J., Quick, J., Simpson, J.T. (2015) A complete bacterial genome assem-
605 bled *de novo* using only Nanopore sequencing data. *Nat. Methods* **12** (8): 733-735.
606 <https://doi.org/10.1038/nmeth.3444>
- 607 [2] Wick, R.R., Judd, L.M., Gorrie, C.L., Holt, K.E. (2017). Completing bacterial genome
608 assemblies with multiplex MinION sequencing, *Microbial Genomics* **3**(10): e000132.
609 <https://doi.org/10.1099/mgen.0.000132>
- 610 [3] Doyle, L.E., Williams, R.B.H., Rice, S.A., Marsili, E., Lauro, F.M. (2018). Draft
611 genome sequence of *Enterobacter* sp. Strain EA-1, an electrochemically active mi-
612 croorganism isolated from tropical sediment, *Genome Announcements* **6**(9): e00111-18.
613 <https://doi.org/10.1128/genomeA.00111-18>
- 614 [4] Daebeler, A., Herbold C.W., Vierheilig, J., Sedlacek C.J., Pjevac, P., Al-
615 bertsen, M., Kirkegaard, R.H., de la Torre, J.R., Daims, H., Wagner, M.
616 (2018). Cultivation and genomic analysis of “*Candidatus Nitrosocaldus islandi-*
617 *cus*”, an obligately thermophilic, ammonia-oxidizing Thaumarchaeon from a hot
618 spring biofilm in Graendalur Valley, Iceland. *Frontiers in Microbiology* **9**: 193.
619 <https://doi.org/10.3389/fmicb.2018.00193>
- 620 [5] Frank, J., Lückner, S., Vossen, R.H.A.M., Jetten, M.S.M., Hall, R.J., Op den Camp,
621 H.J.M., Anvar, S.Y. (2018). Resolving the complete genome of *Kuenenia stuttgartiensis*
622 from a membrane bioreactor enrichment using Single-Molecule Real-Time sequencing.
623 *Scientific Reports*. **8**(1): 4580. <https://doi.org/10.1038/s41598-018-23053-7>

- 624 [6] Andersen, M.H., McIlroy, S.J., Nierychlo, M., Nielsen, P.H., Albertsen,
625 M. (2018). Genomic insights into *Candidatus* Amarolinea aalborgensis gen.
626 nov., sp. nov., associated with settleability problems in wastewater treatment
627 plants, *Systematic and Applied Microbiology*, available online 16 August 2018
628 <https://doi.org/10.1016/j.syapm.2018.08.001>
- 629 [7] Driscoll, C.B., Otten T.G., Brown, N.B., Dreher, T.W. (2017). Towards long-read
630 metagenomics: complete assembly of three novel genomes from bacteria dependent on
631 a diazotrophic cyanobacterium in a freshwater lake co-culture, *Standards in Genomic
632 Sciences*. **12**: 9. <https://dx.doi.org/10.1186/s40793-017-0224-8>
- 633 [8] Slaby, B.M., Hackl, T., Horn, H., Bayer, K., Hentschel, U. (2017). Metagenomic binning
634 of a marine sponge microbiome reveals unity in defense but metabolic specialization,
635 *ISME Journal*, **11**: 2465–2478. <https://doi.org/10.1038/ismej.2017.101>
- 636 [9] Frank, J.A., Pan, Y., Tooming-Klunderud, A., Eijsink, V.G.H., McHardy, A.C.,
637 Nederbragt, A.J. (2016). Improved metagenome assemblies and taxonomic binning
638 using long-read circular consensus sequence data, *Scientific Reports*, **6**: 25373.
639 <https://doi.org/10.1038/srep25373>
- 640 [10] Sevim, V., Lee, J., Egan, R., Clum, A., Hundley, H., Lee, J., Everroad, R.C., De-
641 tweiler, A.M., Bebout, B.M., Pett-Ridge, J., Gker, M., Murray, A.E., Lindemann,
642 S.R., Klenk, H.P., O'Malley, R., Zane, M., Cheng, J.F., Copeland, A., Daum, C.,
643 Singer, E., Woyke, .T (2019). Shotgun metagenome data of a defined mock community
644 using Oxford Nanopore, PacBio and Illumina technologies. *Scientific Data* **6**(1): 285.
645 <https://doi.org/10.1038/s41597-019-0287-z>.
- 646 [11] Brown, B.L., Watson, M., Minot, S.S., Rivera, M.C. Franklin, R.B. (2017). Min-
647 ION nanopore sequencing of environmental metagenomes: a synthetic approach. *Giga-
648 Science* **6**: 1–10. <https://doi.org/10.1093/gigascience/gix007>
- 649 [12] Nanopore GridION and PromethION Mock Microbial Community Data Community
650 Release, Release 2 (2018-10-17). <https://github.com/LomanLab/mockcommunity>
- 651 [13] Huson, D.H., Albrecht, B., Bagci, C., Bessarab, I., Gorska, A., Jolic, D., Williams,
652 R.B.H (2018). MEGAN-LR: New algorithms allow accurate binning and easy in-
653 teractive exploration of metagenomic long reads and contigs, *Biology Direct* **13**: 6.
654 <https://doi.org/10.1186/s13062-018-0208-7>
- 655 [14] Dilthey, A.T., Jain, C., Koren, S. et al. (2019). Strain-level metagenomic assignment
656 and compositional estimation for long reads with MetaMaps. *Nature Communications*
657 **10**: 3066. <https://doi.org/10.1038/s41467-019-10934-2>.

- 658 [15] Laczny, C.C., Kiefer, C., Galata, V., Fehlmann, T., Backes, C., Keller,
659 A. (2017). BusyBee Web: metagenomic data analysis by bootstrapped su-
660 pervised binning and annotation. *Nucleic Acids Res.* **45** (W1): W171–W179.
661 <https://doi.org/10.1093/nar/gkx348>
- 662 [16] Arumugam, K., Bağcı, C., Bessarab, I., Beier, S., Buchfink, B., Górska, A.,
663 Qiu, G., Huson, D.H., Williams, R.B.H. (2019). Annotated bacterial chromo-
664 somes from frame-shift-corrected long-read metagenomic data, *Microbiome* **7**(1): 61.
665 <https://doi.org/10.1186/s40168-019-0665-y>
- 666 [17] Nicholls, S.M., Quick, J.C., Tang, S.Q., Loman, N.J. (2019). Ultra-deep, long-read
667 nanopore sequencing of mock microbial community standards, *GigaScience*, **8** (5):
668 giz043 <https://doi.org/10.1093/gigascience/giz043>
- 669 [18] Somerville, V., Lutz, S., Schmid, M., Frei, D., Moser, A., Irmeler, S., Frey, J.E., Ahrens,
670 C.H. (2019). Long read-based de novo assembly of low complex metagenome samples
671 results in finished genomes and reveals insights into strain diversity and an active phage
672 system, *BMC Microbiology*, **19**: 143. <https://doi.org/10.1186/s12866-019-1500-0>
- 673 [19] Bertrand, D., Shaw J., Kalathiappan, M., Ng, A.H.Q., Muthiah, S., Li, C.H.,
674 Dvornicic, M., Paliska Soldo, J., Koh, J.Y., Ng, O.T., Barkham, T., Young, B.,
675 Marimuthu, K., Chng, K.R., Sikic, M., Nagarajan, N. (2019). Hybrid metage-
676 nomic assembly enables high-resolution analysis of resistance determinants and mo-
677 bile elements in human microbiomes, *Nature Biotechnology*, **37** (8): 937–944.
678 <https://doi.org/10.1038/s41587-019-0191-2>
- 679 [20] Stewart, R.D., Auffret, M.D., Warr, A., Walker, A.W., Roehe, R., Watson, M.
680 (2019). Compendium of 4,941 rumen metagenome-assembled genomes for rumen
681 microbiome biology and enzyme discovery, *Nature Biotechnology*, **37** (8):953–961.
682 <https://doi.org/10.1038/s41587-019-0202-3>
- 683 [21] Moss, E.L., Maghini, D.G., Bhatt, A.S. (2020). Complete, closed bacte-
684 rial genomes from microbiomes using nanopore sequencing. *Nature Biotechnology*.
685 <https://doi.org/10.1038/s41587-020-0422-6>
- 686 [22] Arumugam, K., Bessarab, I., Liu, X.H., Natarajan, G., Drautz-Moses, D.I., Wuertz,
687 S., Lauro, F.M., Law, Y.Y., Huson, D.H., Williams, R.B.H. (2018). Improving recovery
688 of member genomes from enrichment reactor microbial communities using MinION-
689 based long read metagenomics, *bioRxiv* 465328; <https://doi.org/10.1101/465328>
- 690 [23] Schlegel, H.G., Jannasch, H.W. (1967). Enrichment cultures. *Annual Review of Mi-*
691 *crobiology* **21**: 49–70. <https://doi.org/10.1146/annurev.mi.21.100167.000405>

- 692 [24] Strous, M., Kuenen, J.G., Fuerst, J.A., Wagner, M., Jetten, M.S. (2002). The anam-
693 mox case—a new experimental manifesto for microbiological eco-physiology. *Antonie*
694 *Van Leeuwenhoek*. **81**(1–4):693–702. <https://doi.org/10.1023/a:1020590413079>.
- 695 [25] Wu, L, Ning, D., Zhang, B., Li, Y., Zhang, P., Shan, X., Zhang, Q., Brown, M.R.,
696 Li, Z., Van Nostrand, J.D., Ling, F., Xiao, N., Zhang, Y., Vierheilig, J., Wells,
697 G.F., Yang, Y., Deng, Y., Tu, Q., Wang, A., Global Water Microbiome Consor-
698 tium, Zhang, T., He, Z., Keller, J., Nielsen, P.H., Alvarez, P.J.J., Criddle, C.S.,
699 Wagner, M., Tiedje, J.M., He, Q., Curtis, T.P., Stahl, D.A., Alvarez-Cohen, L.,
700 Rittmann, B.E., Wen, X., Zhou, J. (2019) Global diversity and biogeography of bacte-
701 rial communities in wastewater treatment plants. *Nature Microbiology*. **4**(7): 1183–1195.
702 <https://doi.org/10.1038/s41564-019-0426-5>.
- 703 [26] Quince, C., Walker, A.W., Simpson, J.T., Loman, N.J., Segata, N. (2017) Shot-
704 gun metagenomics, from sampling to analysis, *Nature Biotechnology* **35**: 833–84.
705 <https://doi.org/10.1038/nbt.3935>
- 706 [27] Bowers, R.M., Kyrpides, N.C. Stepanauskas, R., Harmon-Smith, M., Doud, D. et
707 al. (2017). Minimum information about a single amplified genome (MISAG) and a
708 metagenome-assembled genome (MIMAG) of bacteria and archaea, *Nature Biotech-*
709 *nology* **35**: 725–731. <https://doi.org/10.1038/nbt.3893>
- 710 [28] Tillett, D., Neilan, B.A. (2000). Xanthogenate nucleic acid isolation from cul-
711 tured and environmental cyanobacteria. *Journal of Phycology* **36**(1): 251–258.
712 <https://doi.org/10.1046/j.1529-8817.2000.99079.x>
- 713 [29] Porechop: <https://github.com/rrwick/Porechop>
- 714 [30] Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., Phillippy,
715 A.M. (2017). Canu: scalable and accurate long-read assembly via adaptive
716 *k*-mer weighting and repeat separation, *Genome Research*, **27**(5): 722–736.
717 <https://doi.org/10.1101/gr.215087.116>
- 718 [31] Wick, R.R., Judd, L.M., Gorrie, C.L. Holt, K.E. (2017). Unicycler: Resolving bacterial
719 genome assemblies from short and long sequencing reads, *PLoS Computational Biology*
720 **13**(6): e1005595. <https://doi.org/10.1371/journal.pcbi.1005595>
- 721 [32] Kolmogorov, M., Rayko, M., Yuan, J., Pevnikov, E., Pevzner, P. (2019). metaFlye:
722 scalable long-read metagenome assembly using repeat graphs, *bioRxiv* 637637;
723 <https://doi.org/10.1101/637637>.
- 724 [33] Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
725 **34** (18): 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>

- 726 [34] Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
727 Abecasis, G., Durbin, R., 1000 Genome Project Data Processing Subgroup (2009) The
728 Sequence Alignment/Map format and SAMtools, *Bioinformatics*. **25**(16): 2078–2079.
729 <https://doi.org/10.1093/bioinformatics/btp352>
- 730 [35] Buchfink, B., Xie, C., Huson, D.H. (2015). Fast and sensitive protein alignment using
731 DIAMOND, *Nature Methods* **12**(1): 59–60. <https://doi.org/10.1038/nmeth.3176>
- 732 [36] O’Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D *et al.* (2016).
733 Reference sequence (RefSeq) database at NCBI: current status, taxonomic ex-
734 pansion, and functional annotation. *Nucleic Acids Res.* **44**(D1): D733–745.
735 <https://doi.org/10.1093/nar/gkv1189>.
- 736 [37] Huson, D.H., Beier, S., Flade, I., Grska, A., El-Hadidi, M., Mitra, S., Ruscheweyh,
737 H.J., Tappu, R. (2016) MEGAN Community Edition – Interactive Exploration and
738 Analysis of Large-Scale Microbiome Sequencing Data. *PLoS Computational Biology*
739 **12**(6): e1004957. <https://doi.org/10.1371/journal.pcbi.1004957>.
- 740 [38] Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., Tyson, G.W.
741 (2015). CheckM: assessing the quality of microbial genomes recovered from
742 isolates, single cells, and metagenomes *Genome Research*, **25**, 1043–1055.
743 <https://doi.org/10.1101/gr.186072.114>.
- 744 [39] Seemann, T. (2014). Prokka: Rapid prokaryotic genome annotation, *Bioinformatics*
745 **30**(14): 2068–2069. <https://doi.org/10.1093/bioinformatics/btu153>
- 746 [40] Olm, M.R., Brown, C.T., Brooks, B., Banfield, J.F. (2018). dRep: a tool
747 for fast and accurate genomic comparisons that enables improved genome re-
748 covery from metagenomes through de-replication, *ISME J.* **11**(12): 2864–2868.
749 <https://doi.org/10.1038/ismej.2017.126>.
- 750 [41] Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. (2019). GTDB-Tk: a toolkit
751 to classify genomes with the Genome Taxonomy Database *Bioinformatics*; btz848.
752 <https://doi.org/10.1093/bioinformatics/btz848>
- 753 [42] Pruesse, E., Quast, C., Knittel, K., Fuchs, B.M., Ludwig, W., Peplies, J., Glöckner,
754 F.O. (2007). SILVA: a comprehensive online resource for quality checked and aligned
755 ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research* **35**: 7188–
756 7196. <https://doi.org/10.1093/nar/gkm864>
- 757 [43] Pruesse, E., Peplies, J., Glöckner, F.O. (2012) SINA: accurate high-throughput mul-
758 tiple sequence alignment of ribosomal RNA genes. *Bioinformatics* **28**: 1823–1829.
759 <https://doi.org/10.1093/bioinformatics/bts252>

- 760 [44] Martin, M. (2011). Cutadapt removes adapter sequences from
761 high-throughput sequencing reads. *EMBnet.journal*, **17**(1): 10–12.
762 <https://doi.org/10.14806/ej.17.1.200>
- 763 [45] Nurk, S., Meleshko, D., Korobeynikov, A., Pevzner, P.A. (2017). metaS-
764 PAdes: a new versatile metagenomic assembler. *Genome Research* **27**(5): 824–834.
765 <https://doi.org/10.1101/gr.213959.116>
- 766 [46] Kang, D.D., Froula, J., Egan, R., Wang, Z. (2015). MetaBAT, an efficient tool for
767 accurately reconstructing single genomes from complex microbial communities, *PeerJ*,
768 **3**, e1165. <https://doi.org/10.7717/peerj.1165>
- 769 [47] Edgar, R.C. (2017). SEARCH_16S: A new algorithm for iden-
770 tifying 16S ribosomal RNA genes in contigs and chromosomes.
771 <http://biorxiv.org/content/early/2017/04/04/124131>
- 772 [48] Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J. (1990). Ba-
773 sic local alignment search tool, *Journal of Molecular Biology* **215**: 403–410.
774 [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- 775 [49] Watson, M., Warr, A. (2019). Errors in long-read assemblies can crit-
776 ically affect protein prediction, *Nature Biotechnology* **37**: 124–126.
777 <https://doi.org/10.1038/s41587-018-0004-z>
- 778 [50] Medaka. <https://github.com/nanoporetech/medaka>
- 779 [51] Vaser, R., Sović I., Nagarajan, N., Šikić, M. (2017). Fast and accurate de novo
780 genome assembly from long uncorrected reads, *Genome Research* **27**(5): 737–746.
781 <https://doi.org/10.1101/gr.214270.116>
- 782 [52] Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz,
783 G., Mesirov, J.P. (2011). Integrative Genomics Viewer. *Nature Biotechnology* **29**: 24–
784 26. <https://doi.org/10.1038/nbt.1754>
- 785 [53] BCFtools. <https://github.com/samtools/bcftools/>
- 786 [54] Skennerton, C.T., Barr, J.J., Slater, F.R., Bond, P.L., Tyson, G.W. (2015). Expanding
787 our view of genomic diversity in *Candidatus* Accumulibacter clades, *Environmental*
788 *Microbiology* **17**(5): 1574–1585. <https://doi.org/10.1111/1462-2920.12582>
- 789 [55] Qiu, G., Liu, X., Saw, N.M.M.T., Law, Y.Y., Zuniga-Montanez, R., Thi, S.S., Ngoc
790 Nguyen, T.Q., Nielsen, P.H., Williams, R.B.H., Wuertz, S. (2019) Metabolic traits of
791 *Candidatus* Accumulibacter clade IIF Strain SCELSE-1 using amino acids as carbon
792 sources for enhanced biological phosphorus removal. *Environmental Science and Tech-*
793 *nology* **54**(4): 2448–2458. <https://doi.org/10.1021/acs.est.9b02901>.

- 794 [56] Stokholm-Bjerregaard, M., McIlroy S.J., Nierychlo, M., Karst, S.M., Al-
795 bertsen, M., Nielsen, P.H. (2017). A critical assessment of the microorgan-
796 isms proposed to be important to enhanced biological phosphorus removal in
797 full-scale wastewater treatment systems, *Frontiers in Microbiology* **8**: 718.
798 <https://doi.org/10.3389/fmicb.2017.00718>
- 799 [57] McIlroy SJ, Onetto CA, McIlroy B, Herbst FA, Dueholm MS, Kirkegaard RH,
800 Fernando E, Karst SM, Nierychlo M, Kristensen JM, Eales KL, Grbin PR, Wim-
801 mer R, Nielsen PH. (2018) Genomic and *in situ* analyses reveal the *Micropru-*
802 *ina* spp. as abundant fermentative glycogen accumulating organisms in enhanced
803 biological phosphorus removal systems. *Frontiers in Microbiology*. **23**(9): 1004.
804 <https://doi.org/10.3389/fmicb.2018.01004>
- 805 [58] Shintani T, Liu WT, Hanada S, Kamagata Y, Miyaoka S, Suzuki T, Naka-
806 mura K. (2000) *Micropruina glycogenica* gen. nov., sp. nov., a new Gram-
807 positive glycogen-accumulating bacterium isolated from activated sludge. *Inter-*
808 *national Journal of Systemic and Evolutionary Microbiology* **50**(1): 201–207.
809 <https://doi.org/10.1099/00207713-50-1-201>
- 810 [59] Onetto, C.A., Grbin, P.R., McIlroy, S.J., Eales, K.L. (2019) Genomic insights into
811 the metabolism of 'Candidatus Defluviicoccus seviourii', a member of Defluviicoccus
812 cluster III abundant in industrial activated sludge. *FEMS Microbiology Ecology*. **95**(2),
813 *fy231* <https://doi.org/10.1093/femsec/fiy231>
- 814 [60] Nielsen, P.H., De Muro, M.A., Nielsen, J.L.. (2000) Studies on the in situ physiology
815 of *Thiothrix* spp. present in activated sludge. *Environmental Microbiology* **2**: 389–398.
816 <https://doi.org/10.1046/j.1462-2920.2000.00120.x>
- 817 [61] Rossetti, S., Blackall, L.L., Levantesi, C., Uccelletti, D., Tandoi, V.
818 (2003). Phylogenetic and physiological characterization of a heterotrophic,
819 chemolithoautotrophic *Thiothrix* strain isolated from activated sludge. *Inter-*
820 *national Journal of Systematic and Evolutionary Microbiology* **53**: 1271–1276.
821 <https://doi.org/10.1099/ijs.0.02647-0>
- 822 [62] Speth, D.R., in 't Zandt, M.H., Guerrero-Cruz, S., Dutilh, B.E., Jet-
823 ten, M.S.M. (2016) Genome-based microbial ecology of anammox granules in
824 a full-scale wastewater treatment system. *Nature Communications* **7**: 11172.
825 <https://doi.org/10.1038/ncomms11172>
- 826 [63] Parks, D.H., Rinke, C., Chuvochina, M., Chaumeil, P.A., Woodcroft, B.J., Evans,
827 P.N., Hugenholtz, P., Tyson, G.W. (2017). Recovery of nearly 8,000 metagenome-
828 assembled genomes substantially expands the tree of life, *Nature Microbiology* **2**(11):
829 1533–1542. <https://doi.org/10.1038/s41564-017-0012-7>

- 830 [64] Kantor, R.S., van Zyl, A.W., van Hille, R.P., Thomas, B.C., Harrison, S.T., Banfield,
831 J.F. (2015) Bioreactor microbial ecosystems for thiocyanate and cyanide degradation
832 unravelled with genome-resolved metagenomics. *Environmental Microbiology* **17**(12):
833 4929–4941. <https://doi.org/10.1111/1462-2920.12936>
- 834 [65] Barbosa, V.L., Atkins, S.D., Barbosa, V.P., Burgess, J.E., Stuetz, R.M. (2006). Char-
835 acterization of *Thiobacillus thioparus* isolated from an activated sludge bioreactor used
836 for hydrogen sulfide treatment. *Journal of Applied Microbiology* **101**(6): 1269–1281.
837 <https://doi.org/10.1111/j.1365-2672.2006.03032.x>
- 838 [66] Crapart, S., Fardeau, M.L., Cayol, J.L., Thomas, P., Sery, C., Ollivier, B., Combet-
839 Blanc, Y. (2007). *Exiguobacterium profundum* sp. nov., a moderately thermophilic,
840 lactic acid-producing bacterium isolated from a deep-sea hydrothermal vent. *In-*
841 *ternational Journal of Systematic and Evolutionary Microbiology* **57**(2): 287–292.
842 <https://doi.org/10.1099/ijs.0.64639-0>
- 843 [67] Mohan Kulshreshtha, N., Kumar, R., Begum, Z., Shivaji, S., Kumar, A. (2013). *Ex-*
844 *iguobacterium alkaliphilum* sp. nov. isolated from alkaline wastewater drained sludge of
845 a beverage factory. *International Journal of Systematic and Evolutionary Microbiology*
846 **63**(12): 4374–4379. <https://doi.org/10.1099/ijs.0.039123-0>
- 847 [68] Lee, S.H., Chung, C.W., Yu, Y.J., Rhee, Y.H. (2009) Effect of alkaline protease-
848 producing *Exiguobacterium* sp. YS1 inoculation on the solubilization and bacterial
849 community of waste activated sludge. *Bioresource Technology* **100**(20): 4597–4603.
850 <https://doi.org/10.1016/j.biortech.2009.04.056>
- 851 [69] Collingro, A., Poppert, S., Heinz, E., Schmitz-Esser, S., Essig, A., Schweikert, M.,
852 Wagner, M., Horn, M. (2005). Recovery of an environmental chlamydia strain from
853 activated sludge by co-cultivation with *Acanthamoeba* sp. *Microbiology* **151**: 301–30.
854 <https://doi.org/10.1099/mic.0.27406-0>.
- 855 [70] Madoni, P. (2011) Protozoa in wastewater treatment processes: A minireview, *Italian*
856 *Journal of Zoology* **78**(1): 3–11, <https://doi.org/10.1080/11250000903373797>
- 857 [71] Mao, Y., Yu, K., Xia, Y., Chao, Y., Zhang, T. (2014). Genome reconstruction and gene
858 expression of “*Candidatus Accumulibacter phosphatis*” Clade IB performing biologi-
859 cal phosphorus removal. *Environmental Science and Technology* **48**(17): 10363–10371.
860 <https://doi.org/10.1021/es502642b>
- 861 [72] Ryu, S.H., Park, M., Lee, J.R., Yun, P.Y., Jeon, C.O. (2007). *Brevundi-*
862 *monas aveniformis* sp. nov., a stalked species isolated from activated sludge. *In-*
863 *ternational Journal of Systematic and Evolutionary Microbiology* **57**(7):1561–1565.
864 <https://doi.org/10.1099/ijs.0.64737-0>

- 865 [73] Wang, C., Zhang, M., Cheng, F., Geng, Q (2015). Biodegradation characterization
866 and immobilized strains' potential for quinoline degradation by *Brevundimonas* sp. K4
867 isolated from activated sludge of coking wastewater. *Bioscience, Biotechnology and Bio-*
868 *chemistry* textbf79(1): 164–170. <https://doi.org/10.1080/09168451.2014.952615>
- 869 [74] Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., Begh-
870 ini, F., Manghi, P., Tett, A., Ghensi, P., Collado, M.C., Rice, B.L., DuLong, C.,
871 Morgan, X.C., Golden, C.D., Quince, C., Huttenhower, C., Segata, N. (2019). Ex-
872 tensive unexplored human microbiome diversity revealed by over 150,000 genomes
873 from metagenomes spanning age, geography, and lifestyle. *Cell* **176**(3): 649–662.e20.
874 <https://doi.org/10.1016/j.cell.2019.01.001>
- 875 [75] Stewart, R.D., Auffret, M.D., Warr, A., Wiser, A.H., Press, M.O., Langford, K.W., Li-
876 achko, I., Snelling, T.J., Dewhurst, R.J., Walker, A.W., Roehe, R., Watson, M. (2018).
877 Assembly of 913 microbial genomes from metagenomic sequencing of the cow rumen,
878 *Nature Communications* **9**: 870. <https://doi.org/10.1038/s41467-018-03317-6>.
- 879 [76] Quince, C., Delmont, T.O., Raguideau, S., Alneberg, J., Darling,
880 A.E., Collins, G., Eren, A.M. (2017). DESMAN: a new tool for de
881 novo extraction of strains from metagenomes, *Genome Biol.* **18**(1): 181
882 <https://doi.org/10.1186/s13059-017-1309-9>.
- 883 [77] Delmont, D.O., Quince, C., Shaiber, A., Esen, O.E., Lee, S.T.M., Rappé, M.S., McLel-
884 lan, S.L., Lückner, S., Eren, A.M. (2018). Nitrogen-fixing populations of Planctomycetes
885 and Proteobacteria are abundant in surface ocean metagenomes, *Nature Microbiology*
886 **3**: 804–813. <https://doi.org/10.1038/s41564-018-0176-9>.
- 887 [78] Ji, P., Zhang, Y.M., Wang, J.F., Zhao, F.Q. (2017). MetaSort untangles metagenome
888 assembly by reducing microbial community complexity. *Nature Communications*, **8**,
889 14306. <https://doi.org/10.1038/ncomms14306>.
- 890 [79] Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A.M., Banfield, J.F.
891 (2019). Accurate and Complete Genomes from Metagenomes, *bioRxiv* 808410 doi:
892 <https://doi.org/10.1101/808410>
- 893 [80] Burton, J.N., Liachko, I., Dunham, M.J., Shendure, J. (2014). Species-level decon-
894 volution of metagenome assemblies with Hi-C-based contact probability maps. *G3*
895 (*Bethesda*), **4**(7): 1339–46. <https://doi.org/10.1534/g3.114.011825>
- 896 [81] Marbouty, M., Cournac, A., Flot, J.F., Marie-Nelly, H., Mozziconacci, J.,
897 Koszul, R. (2014). Metagenomic chromosome conformation capture (meta3C) un-
898 veils the diversity of chromosome organization in microorganisms. *Elife* **3**: e03318.
899 <https://doi.org/10.7554/eLife.03318>

- 900 [82] DeMaere, M., Darling, A. (2019). bin3C: exploiting Hi-C sequencing data
901 to accurately resolve metagenome-assembled genomes. *Genome Biology* **20**: 46.
902 <https://doi.org/10.1186/s13059-019-1643-1>
- 903 [83] Bishara, A., Moss, E.L., Kolmogorov, M., Parada, A.E., Weng, Z., Sidow, A., Dekas,
904 A.E., Batzoglou, S., Bhatt, A.S. (2018) High-quality genome sequences of uncul-
905 tured microbes by assembly of read clouds. *Nature Biotechnology* **36**: 1067–1075.
906 <https://doi.org/10.1038/nbt.4266>
- 907 [84] Sanders, J.G., Nurk, S., Salido, R.A., Minich, J., Xu, Z.Z., Zhu, Q., Martino,
908 C., Fedarko, M., Arthur, T.D., Chen, F., Boland, B.S., Humphrey, G.C., Bren-
909 nan, C., Sanders, K., Gaffney, J., Jepsen, K., Khosroheidari, M., Green, C.,
910 Liyanage, M., Dang, J.W., Phelan, V.V., Quinn, R.A., Bankevich, A., Chang,
911 J.T., Rana, T.M., Conrad, D.J., Sandborn, W.J., Smarr, L., Dorrestein, P.C.,
912 Pevzner, P.A., Knight, R. (2019) Optimizing sequencing protocols for leaderboard
913 metagenomics by combining long and short reads. *Genome Biology* **20**(1):226.
914 <https://doi.org/10.1186/s13059-019-1834-9>.
- 915 [85] Rang, F.J., Kloosterman, W.P., de Ridder, J. (2018). From squiggle to basepair:
916 computational approaches for improving nanopore sequencing read accuracy. *Genome*
917 *Biology* **19**: 90. <https://doi.org/10.1186/s13059-018-1462-9>

Table 1: Summary statistics for 21 putative genomes recovered in this study

Genome identifier	Length (bp)	#CDS ^a	#rRNA ^a	#tRNA ^a	Completeness ^b	Contamination ^b	Taxonomic annotation ^c
PAO1-tig00000001	5,190,177	5,116	2	53	95.28	1.11	s__Accumulibacter sp005584975
PAO1-tig00000003	4,268,816	5,123	1	36	92.82	0.25	g__OLB11
PAO1-tig00000117	2,656,706	3,153	1	38	97.53	0.15	f__UBA6002
PAO1-tig00026549	4,352,448	4,225	1	46	94.64	0.50	f__2-12-FULL-67-15
PAO1-tig00026557	3,138,394	3,619	2	44	94.19	1.58	f__Parachlamydiaceae
PAO1-tig00026560	4,262,704	4,373	2	37	93.99	0.55	g__ELB16-189
PAO1-tig00198536	3,913,768	3,521	2	38	92.57	1.98	s__OLB8 sp001567405
PAO2-tig00000001 ^d	5,027,886	4,558	2	46	93.85	0.00	g__Accumulibacter
PAO2-tig00000013	3,452,123	3,382	2	46	94.49	0.18	g__Dechloromonas
PAO3A-tig00000003	3,666,458	3,454	1	53	94.56	0.28	g__Micropruina
PAO3A-tig00000024 ^d	2,740,818	2,753	1	44	96.92	0.00	s__Brevundimonas sp002426005
PAO3A-tig00000209	3,302,829	3,190	1	47	95.38	0.00	f__Nocardioideaceae
PAO3A-tig00018026	4,685,957	4,742	1	48	93.47	1.09	c__Anaerolineae
PAO3A-tig00139797 ^d	3,548,924	3,508	2	47	99.04	0.48	s__Thiobacillus sp00189930
PAO3B-tig00000024 ^d	3,282,734	2,937	2	51	98.10	0.23	f__Burkholderiaceae
PAO3B-tig00000027 ^d	3,375,962	3,179	1	53	93.16	0.42	g__Rhodoblastus
PAO4-tig00000001 ^d	3,950,501	3,652	2	58	98.64	2.07	g__Pseudoxanthomonas_A
PAO4-tig00000030 ^d	4,541,730	4,623	2	47	97.10	2.30	g__Thiothrix
PAO4-tig00000046 ^d	3,961,963	3,725	2	51	96.30	0.47	g__Dechloromonas
PAO4-tig00000028	4,261,978	4,404	2	50	94.81	0.00	g__UBA1943
PAO4-tig00000079 ^d	2,921,657	3,197	9	59	94.79	0.00	s__Exiguobacterium profundum

^a As predicted using Prokka (see Methods). The rRNA count denotes presence of three genes encoding 5S, 16S and 23S sequences.

^b Genome quality estimates from CheckM (see Methods).

^c Taxonomic assignments from GTDB-Tk (see Methods).

^d Sequence classified as *circular* by Canu.

Table 2: Influence of sequence procedures on CheckM-derived genome quality statistics

Quality measure	Completeness ^a						Contamination ^a									
	Uncorrected			Medaka ^b			MEGAN			Racon ^c			Multiple ^d			
	MEGAN	Medaka ^b	Racon ^c	MEGAN	Medaka ^b	Racon ^c	MEGAN	Medaka ^b	Racon ^c	MEGAN	Medaka ^b	Racon ^c	MEGAN	Medaka ^b	Racon ^c	Multiple ^d
PAO1-tig00000001	84.28	95.28	93.91	87.88	94.71	0.66	1.11	1.59	1.62	1.11	1.59	1.62	1.11	1.59	1.62	1.11
PAO1-tig00000003	64.34	92.82	84.84	78.21	87.74	1.56	0.25	0.90	0.74	0.25	0.90	0.74	0.25	0.90	0.74	0.74
PAO1-tig00000117	82.46	97.53	94.82	90.99	94.82	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15
PAO1-tig00026549	86.57	94.64	97.13	91.45	96.73	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50
PAO1-tig00026557	75.56	94.19	95.24	86.16	95.58	2.70	1.58	1.58	1.58	1.58	1.58	1.58	1.58	1.58	1.58	1.58
PAO1-tig00026560	71.51	93.99	94.24	86.70	94.52	1.59	0.55	0.00	0.30	0.55	0.00	0.30	0.55	0.00	0.30	0.00
PAO1-tig00198536	74.83	92.57	94.64	88.68	94.39	2.74	1.98	2.48	2.23	1.98	2.48	2.23	1.98	2.48	2.23	2.48
PAO2-tig00000001	70.90	93.85	89.33	83.64	90.88	0.03	0.00	0.03	1.19	0.00	0.03	1.19	0.00	0.03	1.19	0.03
PAO2-tig00000013	45.13	94.49	89.16	86.32	91.52	0.00	0.18	0.59	0.38	0.18	0.59	0.38	0.18	0.59	0.38	0.12
PAO3A-tig00000003	75.75	94.56	93.78	83.48	93.52	0.10	0.28	0.50	0.10	0.28	0.50	0.10	0.28	0.50	0.10	0.50
PAO3A-tig00000024	79.79	96.92	89.46	84.66	89.86	0.81	0.00	0.32	0.32	0.00	0.32	0.32	0.00	0.32	0.32	0.32
PAO3A-tig00000209	87.00	95.38	96.80	94.60	97.32	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
PAO3A-tig00018026	41.99	93.47	81.67	66.02	77.17	0.28	1.09	1.09	1.27	1.09	1.09	1.27	1.09	1.09	1.27	1.37
PAO3A-tig00139797	91.56	99.04	97.05	97.02	97.05	0.48	0.48	0.55	0.48	0.48	0.55	0.48	0.48	0.55	0.48	0.48
PAO3B-tig00000024	89.02	98.10	93.63	90.53	91.71	0.23	0.23	0.23	0.00	0.23	0.23	0.00	0.23	0.23	0.00	0.39
PAO3B-tig00000027	73.78	93.16	85.11	81.80	87.97	0.42	0.42	0.63	0.16	0.42	0.63	0.16	0.42	0.63	0.16	0.73
PAO4-tig00000001	93.04	98.64	97.97	96.58	97.95	1.90	2.07	2.07	2.06	2.07	2.07	2.06	2.07	2.07	2.06	2.06
PAO4-tig00000030	76.77	97.10	91.63	85.14	91.87	1.30	2.30	1.52	1.63	2.30	1.52	1.63	2.30	1.52	1.63	1.11
PAO4-tig00000046	91.85	96.30	97.64	94.90	97.40	0.47	0.47	0.47	0.50	0.47	0.47	0.50	0.47	0.47	0.50	0.47
PAO4-tig00000079	80.26	94.79	89.66	89.85	90.87	0.00	0.00	0.66	0.72	0.00	0.66	0.72	0.00	0.66	0.72	0.66
PAO4-tig000000228	70.80	94.81	77.98	77.13	78.18	0.02	0.00	0.08	0.00	0.00	0.08	0.00	0.00	0.08	0.00	0.00

Genomes shown in bold show CheckM completeness >90% and contamination < 5%.

^a Genome completeness and contamination estimates obtained from *CheckM* (see Materials and Methods).

^b Medaka: Single round of correction applied to each uncorrected genome sequence.

^c Racon: Single round of correction applied to each uncorrected genome sequence, run with default parameters for all datasets.

^d Multiple: four sequential applications of Racon (run with default parameters using long read data) and then one application of Medaka (run with default parameters; the following models were used -m r941_min_high_g303 for PAO1, PAO2, PAO3A and PAO3B and -m r941_min_high_g330 for PAO4)

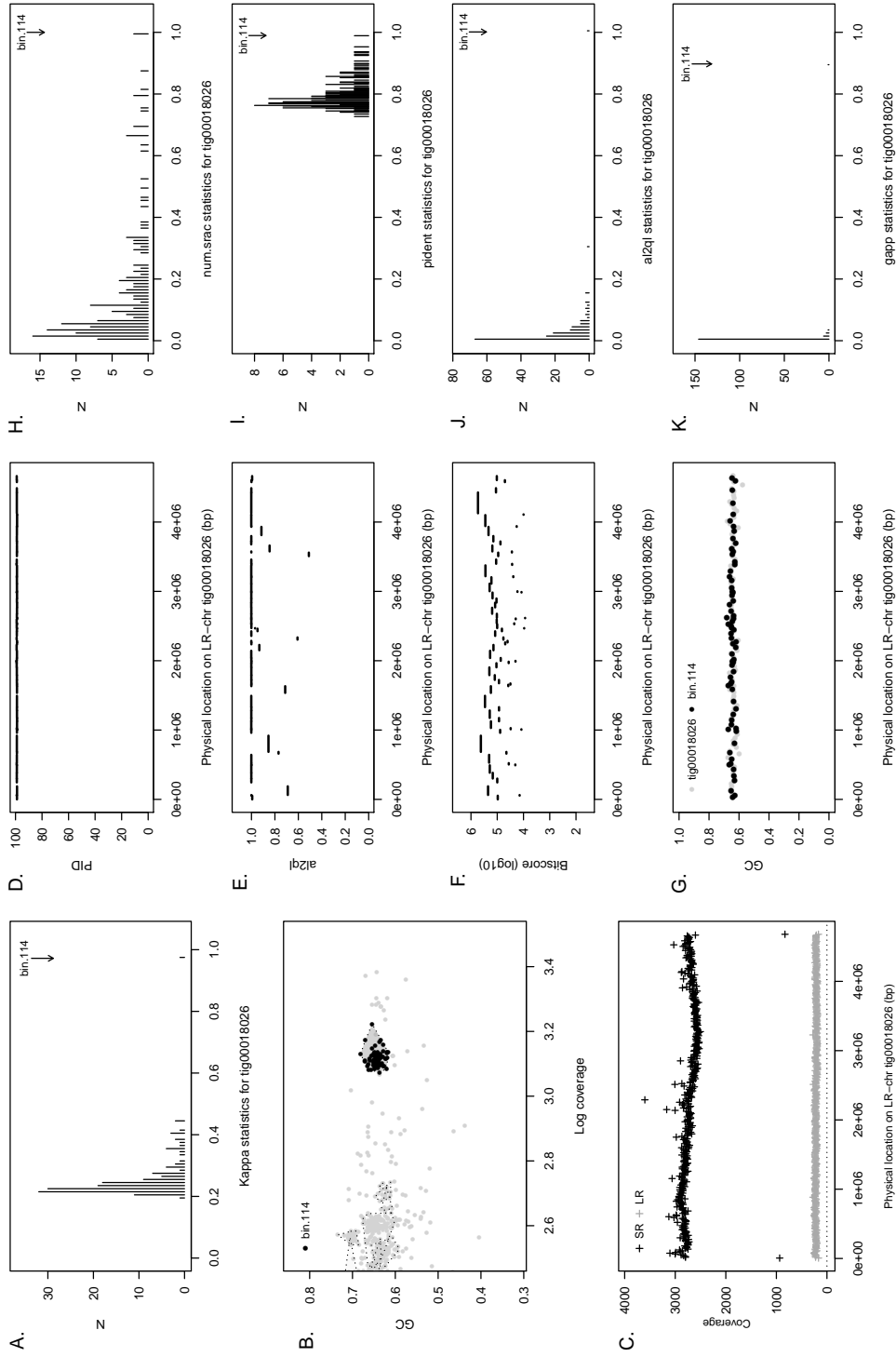


Figure 1: Summary of concordance statistic analysis for an LR-chr (tig00018026) from the PAO3A reactor community (annotated to class *Anaerolineaceae* showing close relationship to a short read metagenome assembled from the same reactor community (bin 114)). (A): Distribution of κ -scores for tig00018026 against 242 bins recovered from the corresponding short read assembly. Bin 114 has the highest κ at 0.97; (B): coverage-GC plot for the short read assembly, with bin 114 highlighted (closed black circles and dark grey convex hull; other bins highlighted by light grey convex hulls); (C): short read (SR, black crosses) and long read (LR, grey crosses) coverage profiles across tig00018026. (D-F): BLASTN statistics for alignments of short read contigs (bin 114) against tig00018026. Horizontal segments show alignment position on LR-chr and height of segment is value of corresponding statistic (y -axis) namely percent identity (PID) (D), the ratio of alignment length to query length (al2ql)(D) and \log_{10} -bitscore (E). (F): GC content as a function of position on tig00018026 (grey closed circles, computing in adjacent windows of length 46700 bp) and for aligned short read contigs (black closed circles); (G-K): distribution of four component statistics of κ (see **Methods**), with the position of the top scoring short read bin highlighted. (G): proportion of short read contigs in bin aligned to LR-chr (p_{srac}); (H): mean percent identity (\widehat{pid}). (I): mean ratio of alignment length to query length $\widehat{al2ql}$ and (K): proportion of the long read contig that is covered by an alignment (p_{aln}).

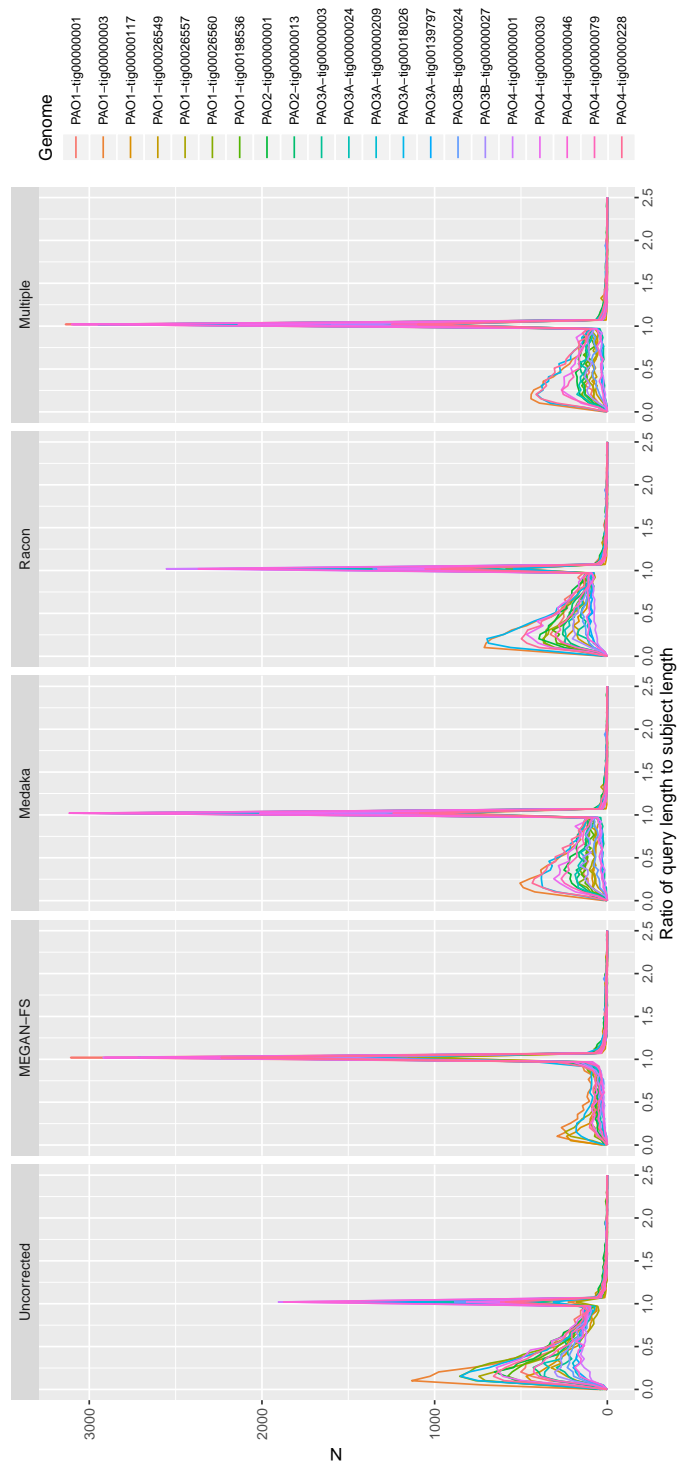


Figure 2: Density estimates for the length ratio statistics, computed from the length of predicted genes in long read assemblies (query) and length of their best hit counterparts in cognate short assemblies (subjects), and categorised by type of sequence correction employed (from left to right, raw assembled sequence [uncorrected], frame-shift correction using MEGAN-LR, sequence correction using Medaka, sequence correction using Racon and application of the multiple procedure approach. Results from individual recovered genomes are highlighted by colour, and x -axis truncated at 2.5 units. A version with a log-scale on the vertical axis is provided in **Supplementary Figure 35**

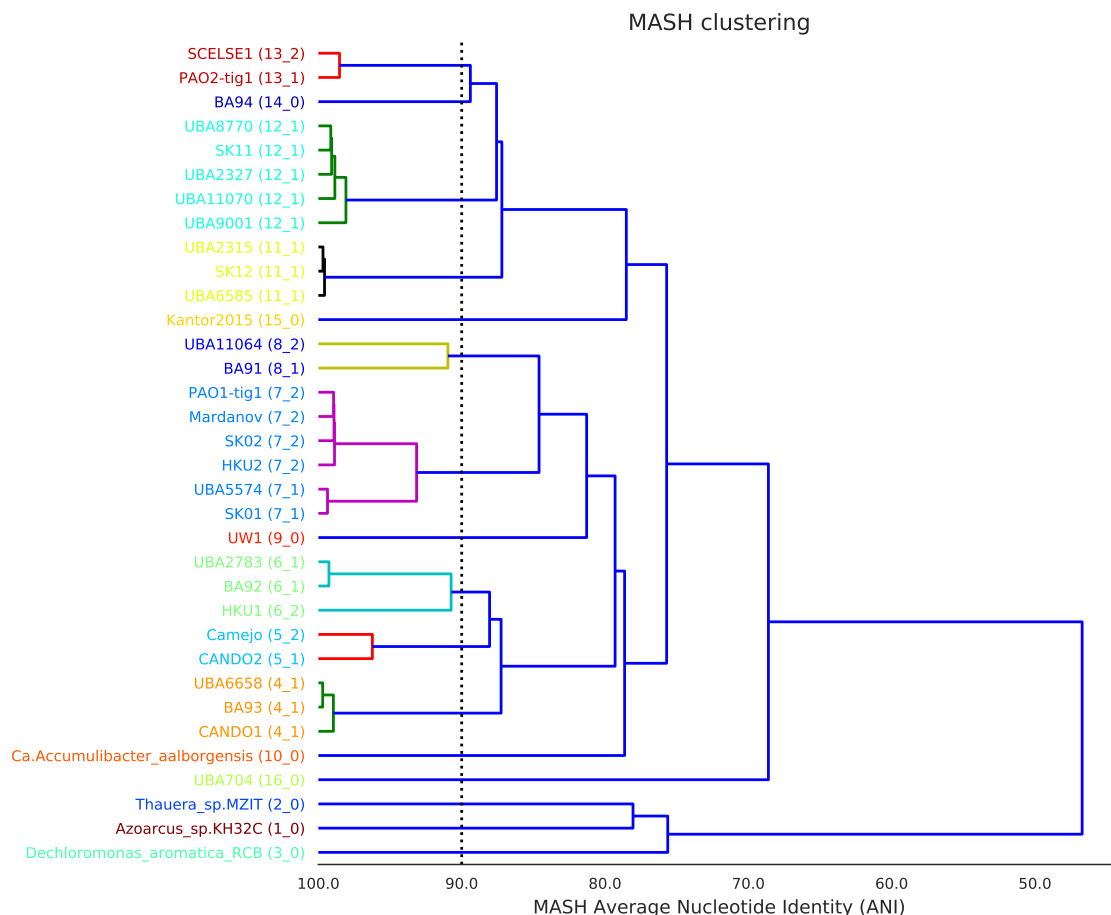


Figure 3: Dendrogram generated from MASH distances between draft genomes of *Candidatus* *Accumulibacter*, including two genomes recovered in the present study. Genomes from genera *Thauera*, *Azoarcus* and *Dechloromonas* were used as an outgroup. Underscore separated number in brackets refers to dRep secondary cluster assignments (two genomes are in the same secondary cluster if their ANImf ≥ 99). Note the structure of the tree recapitulates previously defined clade associations (*Clade IIF*: BA94, SK11, SK12; *Clade IIC*: BA91, SK02, SK01; *Clade I*: BA92 and BA93. With UW1 being a singleton for *Clade IIA*). Genome references as follows, from top of tree: *SCElse1* (GCA_005524045.1); *BA94* (GCA_000585095.1); *UBA2327* (GCA_002345025.1); *SK11* (GCA_000584995.1); *UBA8770* (GCA_003487685.1); *UBA11070* (GCA_003535635.1); *UBA9001* (GCA_003542235.1); *UBA2315* (GCA_002345285.1); *SK12* (GCA_000585015.1); *UBA6585* (GCA_003535635.1); *Banfield* (GCA_001897745.1); *UBA11064* (GCA_003538495.1); *BA91* (GCA_000585035.2); *Mardanov* (GCA_005889575.1); *SK02* (GCA_000584975.2); *HKU2* (GCA_000987395.1); *UBA5574* (GCA_002425405.1); *SK01* (GCA_000584955.2); *UW1* (GCA_000024165.1); *UBA2783* (GCA_002352265.1); *BA92* (GCA_000585055.1); *HKU1* (GCA_000987445.1); *CANDO2* (GCA_009467885.1); *Camejo* (GCA_003332265.1); *UBA6658* (GCA_002455435.1); *BA93* (GCA_000585075.1); *CANDO1* (GCA_009467855.1); *Ca. Accumulibacter aalborgensis* (GCA_900089955.1); *UBA704* (GCA_002304785.1)