

1     **Shallow MinION sequencing to assist *de novo* assembly**  
2             **of the *Streptococcus agalactiae* genome**

3  
4     Tamara Hernandez-Beeftink<sup>1</sup>, Hector Rodriguez-Perez<sup>1</sup>, Ana Díaz-de Usera<sup>2</sup>, Rafaela  
5             Gonzalez-Montelongo<sup>2</sup>, José M. Lorenzo-Salazar<sup>2</sup>, Fabián Lorenzo-Díaz<sup>1,3</sup>, Carlos  
6                             Flores<sup>1,2,4\*</sup>

7  
8     <sup>1</sup>Research Unit, Hospital Universitario Ntra. Sra. de Candelaria, Universidad de La Laguna, Santa Cruz  
9     de Tenerife, Spain; <sup>2</sup>Genomics Division, Instituto Tecnológico y de Energías Renovables (ITER), Santa  
10    Cruz de Tenerife, Spain; <sup>3</sup>Departamento de Bioquímica, Microbiología, Biología Celular y Genética,  
11    Universidad de La Laguna, Santa Cruz de Tenerife, Spain; <sup>4</sup>CIBER de Enfermedades Respiratorias,  
12    Instituto de Salud Carlos III, Madrid, Spain.

13  
14  
15  
16    \*Corresponding author

17    E-mail: [cflores@ull.edu.es](mailto:cflores@ull.edu.es) (CF)

18

19

20

21

22

23

24

25

## 26 **Abstract**

27 Despite the reduced read length, the so-called Next-Generation Sequencing (NGS) of  
28 second-generation has allowed rapid and complete genome characterization of many  
29 species. MinION (Oxford Nanopore Technologies), a portable third-generation NGS  
30 device, enables sequencing of long DNA fragments at low cost. Here we used a low-  
31 coverage MinION sequencing in combination with short-read NGS to improve genome  
32 assembly. We tested this possibility by using MinION R9.0 with Rapid 1D kit and  
33 MiSeq with >300X paired-end 300 bp reads (Illumina, Inc.) for the genome assembly of  
34 a *Streptococcus agalactiae* clinical isolate (2.2 Mb). With as few as 1,171 MinION  
35 reads that covered the genome at 2.4X (the longest read being 186 Kb long), the hybrid  
36 assembly combining MinION and Illumina reads increased the N50 by 4.9-fold  
37 compared to the assembly using Illumina data alone. Almost 50% of the genome was  
38 represented into a single contig (1.02 Mb). Besides, this allowed the full reconstruction  
39 of mobile elements, including a plasmid, and improved gene annotation. Taken  
40 together, our results support that shallow MinION sequencing combined with high-  
41 throughput second-generation NGS constitutes a cost-efficient strategy for the assembly  
42 of whole genomes.

43 **Key words:** Hybrid assembly; Nanopore; artefactual reads

44

## 45 **Background**

46           Next-Generation Sequencing (NGS) technologies have allowed sequencing and  
47 analysis Group B *Streptococcus* (GBS) genomes from different environments and  
48 organisms[1–6]. The presence of prophages and mobile elements in GBS was  
49 associated with the bacterial adaptation and ability to cause infections, likely enhancing  
50 their diversity and spread[7–10]. In fact, the estimates suggest that roughly 50% of  
51 *Streptococcus agalactiae* genome is grouped into islands of pathogenicity, where  
52 virulence genes and mobile elements are present[11]. *S. agalactiae* genome shows a  
53 complex organization and rearrangements, where the presence of a conserved backbone  
54 and 69 variable regions have been observed[12]. These observations justify the  
55 necessity of using *de novo* assembled genomes for better comparative bacterial studies  
56 and adequate tracking of mobile elements. However, the assembly of bacterial genomes  
57 based on short-read NGS technologies (so called second-generation) usually produce  
58 discontinuous sequences[13,14], often requiring the assistance from algorithms,  
59 alternative library approaches or sequencing technologies to provide sequence  
60 continuity.

61

62           Third-generation NGS technologies allow obtaining longer DNA reads from  
63 PCR-free DNA libraries in real time, constituting a solution for resolving repeats[15]  
64 and obtaining complete assembled genomes. In this context, the portability and reduced  
65 costs of sequencing with the MinION device (Oxford Nanopore Technologies) has  
66 facilitated the adoption of this technology in many distinct applications[16]. Besides its  
67 utility for the analysis of structural variants[17,18], SNP determination[19], cytosine  
68 methylation[20], transcriptomics[21], and in-field experiments in extreme  
69 environments[22], MinION has widely used to sequence small whole genomes from

70 virus (10.8 kb)[23], mitochondrial DNA (16 kb)[24], and bacteria (4.6 Mb)[25],  
71 allowing their assembly into single contigs when depth of coverages were above 20-  
72 40X[26]. Nevertheless, this technology continues to be challenged by its high error rate  
73 (>10%)[27], particularly in homopolymeric regions[28], and the generation of reads that  
74 fail to align against the target[29,30]. Because of this, hybrid assemblies leveraging the  
75 per-base quality of high-coverage (>200X) second-generation data with the bridging  
76 capacity of >10X depth-of-coverage third-generation data is nowadays perceived as the  
77 optimal choice for obtaining accurate assemblies[24,27,31–35]. Given the large  
78 variability in sequence throughput offered by this technology well below product  
79 specifications, here we tested the benefit of shallow MinION sequencing data in  
80 assisting the hybrid genome assembly of a *S. agalactiae* clinical isolate (genome size of  
81 2.2 Mb).

82

## 83 **Methods**

### 84 Bacterial isolate, growing conditions and validation

85 We have previously examined the IMESag-rpsI mobile element distribution in  
86 240 whole genomes of Group B *Streptococcus* isolates. The *S. agalactiae* HRC strain  
87 was subjected to a deeper examination to assist the characterization of the mobile  
88 element in that study[36]. *S. agalactiae* HRC is a serotype V representative strain  
89 belonging to the tetracycline-resistance CC1 lineage Tn916-1[37], which was isolated in  
90 2009 from a 60-year-old woman with abdominal sepsis at the Hospital Ramón y Cajal  
91 (Madrid, Spain)[36].

92

93 The bacterial isolate was cultured in TSB (Tryptic Soy Broth, Becton  
94 Dickinson) liquid medium at 37°C for 24 h. DNA extraction was performed with the

95 GenElute Bacterial Genomic DNA kit (Sigma-Aldrich) and was quantified on a Qubit  
96 3.0 fluorometer using the dsDNA BR assay kit (Thermo Fisher Scientific). Its identity  
97 was first validated by 16S rRNA gene sequencing as follows. DNA was amplified by  
98 PCR with 30 cycles of 95°C for 20 s, 50°C for 30 s, and 72°C for 15 s, using the  
99 HotStarTaq master mix kit (QIAGEN), and the primers 1391R (5'-  
100 GACGGGCGGTGWGTRCA-3') and 27F (5'-AGRGTTYGATYMTGGCTCAG-3').  
101 Amplified material was purified with ExoSAP-IT (Thermo Fisher Scientific), subject to  
102 Sanger sequencing with BigDye Terminator v3.1 Cycle Sequencing kit (Thermo Fisher  
103 Scientific), and products purified with DyeEx 2.0 Spin kit (QIAGEN), following the  
104 manufacturer's recommendations. Sequencing products were resolved and basecalled in  
105 a 3500 Genetic Analyzer (Thermo Fisher Scientific). A simple BLAST[38] search on  
106 the sequences obtained verified that the bacterial isolate had a 99% identity to *S.*  
107 *agalactiae*.

108

#### 109 Illumina sequencing

110 DNA from *S. agalactiae* HRC (1 ng) was used to generate sequencing libraries  
111 with the Nextera XT kit (Illumina, Inc.) following the manufacturer's recommendations.  
112 A MiSeq instrument and a MiSeq Reagent kit V3 (Illumina, Inc.) was used for 300 base  
113 paired-end sequencing along with a spike-in of 20% of the PhiX Control v3 (Illumina,  
114 Inc.). Library construction and the sequencing experiment were performed at the  
115 Instituto Tecnológico y de Energías Renovables (ITER).

116

#### 117 MinION sequencing

118 Rapid Sequencing kit (SQK-RAD001) (Oxford Nanopore Technologies) was  
119 used for library preparation starting from 200 ng of *S. agalactiae* HRC purified DNA.

120 Sequencing with MinION used SpotOn Flow Cell Mk I R9 Version (Oxford Nanopore  
121 Technologies) and followed manufacturer's recommendations, except that the  
122 experiment was left running only for 22 h.

123

#### 124 Bioinformatics and statistical procedures

125 MiSeq Reporter Software v1.18.54 (Illumina, Inc.) was used to convert  
126 intensities to basecalls and generate the FASTQ file. MinION basecalling was  
127 performed with MinKNOW GUI 1.1.21 software (Oxford Nanopore Technologies) and  
128 all records were processed with poRe v0.21[39] to extract the obtained reads in FASTQ  
129 file format. Unicycler v0.4.1[40] was used with default parameters for hybrid *de novo*  
130 assembly of *S. agalactiae* HRC genome, combining Illumina and Nanopore datasets,  
131 which involved multiple rounds of short-read polishing. As a reference for comparisons,  
132 Unicycler was also used to assemble the Illumina data alone. Assembly comparisons  
133 were assessed with QUAST v4.0[41] using the *S. agalactiae* SS1 genome  
134 (NZ\_CP010867.1) as the reference. The assembled contigs were visualized with  
135 Bandage v0.8.1[42]. Annotation of genes and elements in the contigs was done with  
136 IonGAP[43]. Mauve v2.4.0[44] was used to identify the insertion sites of the *S.*  
137 *agalactiae* HRC mobile elements in the assembly against the reference  
138 (NZ\_CP010867.1).

139

140 In order to further explore raw MinION reads, per-read GC content distribution  
141 was first assessed to identify outliers (>3 SD from the mean). The existence of  
142 artefactual reads in the output[29] was evaluated through comparisons against the  
143 reconstructed *S. agalactiae* HRC genome, including the two characteristic elements  
144 identified. BLAST v2.6.0[38] alignments were used to isolate unaligned artefactual

145 reads. Bio.SeqUtils library of Biopython release 1.70 was then used to assess their  
146 differences in length (bp), GC content (GC%) and mean quality (Q) score against the  
147 reads aligning to the *S. agalactiae* HRC assembled genome. A non-parametric Mann-  
148 Whitney U-test was used to assess the significance in the comparisons. To test the  
149 possibility that particular pores were generating more artefacts than expected, as has  
150 been suggested elsewhere[29], unaligned reads were mapped back to the 512 nanopores  
151 of the Flow Cell and significance was tested assuming a Poisson distribution. Finally, to  
152 evaluate their impact in *de novo* assembly of *S. agalactiae* HRC genome, we then  
153 excluded all unaligned reads from the dataset, and used Unicycler under the same  
154 conditions to reassess the assemblies.

155

## 156 **Results**

### 157 Summary of sequencing results

158 Sequencing with Illumina resulted in 2,975,356 paired-end reads of 300 bp with  
159 Phred>30, which translates into 354X depth-of-coverage of *S. agalactiae* genome. In  
160 comparison, the MinION run proceeded with a total of just 276 active pores, providing  
161 1,171 reads with a mean length of 4.5 kb and a Q score average of 10.8, and ranging  
162 from 3.8 to 25.4. While the largest proportion of reads accumulated in the 1.86-9.0 kb  
163 range (the N50 of the output was 9.2 kb) the output included reads as long as 186 kb  
164 (which represents almost 10% of the *S. agalactiae* HRC genome). The total throughput  
165 (5.2 Mb) theoretically translates into a 2.4X depth-of-coverage of the bacterial genome,  
166 which is well below product specifications, despite all initial quality controls of the  
167 platform were passed. Finally, the per-read GC content average of reads was 44.3%  
168 (ranging from 10.3% to 95.9%), noticeably larger than that of the reference sequence  
169 (35.5%; NZ\_CP010867.1).

170 Assembly improvements by using MinION data

171       The Illumina data assembly generated 23 contigs (>500 bp) with a N50 of  
172 147,211 bp, and the largest contig being 525,288 bp long. The length of the assembly  
173 from the Illumina data alone was very close to the *S. agalactiae* reference, accounting  
174 for a genome size of 2,144,732 bp, a GC content of 35.3%, and an aligning length  
175 against the reference of 2,049,610 bp (97.9% of consensus identity) (Table 1). In  
176 comparison, although the hybrid assembly including also combining the Illumina and  
177 MinION reads did not allow to obtain a fully resolved bacterial genome, mostly because  
178 of a major unresolved structure containing repetitive elements, rRNA operon copies and  
179 tRNAs (Fig 1), it reduced the number of contigs to less than a half (8 with >500 bp),  
180 having a N50 of 714,293 bp and the largest contig being 1,018,988 bp long. The hybrid  
181 assembly generated a total length of 2,159,060 bp, with 2,062,825 bp aligning against  
182 the reference (98.5% of consensus identity) (Table 1). Both assemblies also evidenced a  
183 small plasmid (2,491 bp) as an independent replicon, which perfectly aligned with the  
184 *Streptococcus oralis tigurinus* 2426 plasmid pST2426 (ASXA01000016.1), along with  
185 two other mobile elements detected in previously studies[36] (Fig 1). In terms of gene  
186 annotation, the hybrid assembly also outperformed the Illumina only assembly (2,153  
187 vs. 2,135 genes, respectively) when compared to the reference (2,195 genes) (Table 1).  
188 Taken together, the small amount of reads provided by MinION (0.04% of total reads  
189 obtained for the bacterial isolate) allowed to increase the N50 of the genome assembly  
190 by 4.9-fold, to cover almost 50% of the bacterial genome in a single contig, and to  
191 improve the annotation of predicted genes in the assembly.

192

193 Characterization of unaligned MinION reads and impact on the assembly



194 A closer inspection of the MinION run records evidenced the existence of reads  
195 with large regions of limited sequence diversity (end to end) that were largely  
196 compatible with artifacts. Just by exploring the GC content of the MinION output, we  
197 detected 27 reads with large deviations ( $>3$  SD from the mean) corresponding to GC  
198 contents below 7.2% or above 77.8%. The reads involved had a mean size of 716 bp,  
199 but were as large as 4.2 kb. However, in order to better identify and characterize the  
200 reads that were most likely artefactual, we used BLAST at 85% identity threshold to  
201 align all reads against the assembled genome. This allowed to identify 68 unaligned  
202 reads with a mean size of 2.0 kb (but as large as 21.6 kb) and a mean Q score of 6.7 (but  
203 as high as 22.9) (Fig 2). Only three of these unaligned reads (0.25% of total) aligned to  
204 sequences from other species or synthetic constructs according to a BLAST search  
205 (Table S1), most probably corresponding to experimental contaminants. Surprisingly,  
206 MinKNOW software only classified 20 of the 68 as failed based on internal  
207 manufacturer's specifications. Unaligned reads differed significantly from reads that  
208 aligned for the GC content (74.2 vs. 42.5%;  $p= 2.5e^{-37}$ ), Q score (6.7 vs. 11.1;  $p= 2.2e^{-21}$ ),  
209 and mean read length (2.0 vs. 4.6 kb;  $p= 1.3e^{-11}$ ) (Fig 2). Strikingly, from the 276  
210 active pores of the run, unaligned reads were obtained from 45, and two pores in  
211 particular experienced a significant accumulation of them (28.9%,  $p<0.019$ ) (S1 Fig).  
212 Finally, a re-assessment of the hybrid assembly evidenced identical results with or  
213 without excluding unaligned reads, supporting that Unicycler was robust to their  
214 presence in the input. For simplicity, these results are not shown.

215

## 216 Discussion

217 Here we demonstrate that the use of very low amount of MinION reads  
218 combined with short-read sequencing data has the capacity to increase the continuity of

219 assembly contigs and plasmid isolation without manual intervention, having the  
220 potential to fully resolve bacterial genomes. This was illustrated with the assembly of *S.*  
221 *agalactiae* HRC genome based on as few as 2.4X MinION depth-of-coverage data,  
222 paving the way for obtaining cost-effective continuum whole-genome sequences as the  
223 nanopore technology throughput increases. While other species may have alternative  
224 requirements, to the best of our knowledge, this is the first study assessing the benefits  
225 of using MinION reads to assist the complete assembly of microbial genomes with  
226 coverages well below 10X (usually assisted with >14X)[14,23,33,34]. Although we  
227 were unable to fully resolve the assembly, our results are quite reassuring considering  
228 the small number of MinION reads used in the experiment (0.04% of total reads),  
229 allowing to obtain almost the 50% of the bacterial genome in a single contig. Besides,  
230 we were able to detect and fully reconstruct a plasmid, previously attained only with  
231 higher depth-of-coverage and longer MinION reads[45].

232

233 High quality noise reads generated by MinION have been described recently in  
234 the literature[29], consisting on reads that did not map to the target, lab contaminants, or  
235 spike-in sequences. Many of those appeared to be artefacts of the technology that were  
236 associated with specific active pores[29], to which some authors have previously  
237 attributed an insufficient read quality[30]. In this study, MinION generated around 5%  
238 of reads of variable quality that did not map to the target. We found that despite a few of  
239 them had partial identity with bacterial sequences that may indicate minor experimental  
240 contaminations as has been shown by others[15,46], most of them consisted in low-  
241 complexity sequences from end to end. We demonstrate that a number of them were  
242 generated from specific active pores during our experiment. Based on these evidences,  
243 we speculate that most of these reads are artefacts produced during the sequencing

244 process and that they are most likely independent from the basecaller being used. In this  
245 study, we demonstrate that they did not negatively affect the hybrid assembly process in  
246 the conditions assessed. However, their impact in other projected applications (such as  
247 shotgun metagenomics or whole-genome sequencing of eukaryotes) or algorithms other  
248 than Unicycler remain unknown. Solutions to reduce or identify and filter out these  
249 reads will need to be implemented.

250

251 With sufficient MinION yield providing enough depth-of-coverage, the standard  
252 analysis is based on long-read and error-tolerant assemblers such as Canu[47]. Because  
253 of the high error rate of this technology at the moment (10-20%)[27,46,48], and the  
254 problems resolving the homopolymers[28], the assembly process from MinION data  
255 alone typically necessitates of a number of error correction steps to produce high quality  
256 genome sequences[49]. Most importantly, this assembly process is computationally  
257 intensive[49,50]. Theoretically, MinION can produce many thousands of reads and a  
258 few dozen Gb per run[49]. However, in practice MinION outputs are highly variable  
259 translating sometimes in a few thousand reads (3,400 in Benitez-Paez et al[51]) and Mb  
260 scale outputs. In our hands, regular experiments are roughly in the scale of 100 Mb per  
261 run. This is because the number of active pores changes during the experiment (between  
262 105 and 426 in Greninger et al[46]), below theoretical maximum in distinct runs, and  
263 the effective number of useful reads represents a fraction of the output[15,46]. Some  
264 studies have reported that a large proportion of the output (>50%) have failed to align  
265 against the target[30,52]. Besides, sudden crashes of the MinION software  
266 controller[14] or the computer connection failures[53] have been reported during a  
267 sequencing experiment. In addition, unofficial information provided by users in The  
268 Nanopore Community report extremely low occupancy of pores (<10%) during the

269 sequencing experiments with Rapid 1D kits, which have been suggested to be related to  
270 the kit production itself. These numbers are comparable to our observations in this study  
271 and all these issues have obvious consequences in the throughput of a MinION  
272 experiment. Therefore, although not the standard, the limited yield obtained by the  
273 MinION sequencing in this study is not uncommon. In this context, a hybrid assembly  
274 with a very shallow coverage of MinION data constitutes an advantage over the  
275 approaches based on MinION data alone for obtaining cost-effective high-quality  
276 continuous bacterial genomes. Besides being computationally faster, this hybrid option  
277 reduces the burden of generating the thousands of long reads that are necessary to obtain  
278 high quality bacterial genome assemblies with MinION data alone (>20X)[26].

279

## 280 **Conclusions**

281 We demonstrate that the combination of a small proportion of long reads with  
282 high-coverage short-read data constitutes a promising strategy to generate high quality,  
283 highly contiguous assemblies, here allowing to nearly complete the *S. agalactiae* HRC  
284 genome. Besides, we detected MinION reads consisting of low-complexity sequences  
285 compatible with artefacts of the technology but that did not affect the assembly.

286

## 287 **Acknowledgments**

288 Funded by Instituto de Salud Carlos III (PI14/00844; PI17/00610; FI17/00177) and  
289 Ministerio de Ciencia, Innovación y Universidades (RTC-2017-6471-1;  
290 MINECO/AEI/FEDER, UE) co-financed by the European Regional Development  
291 Funds, “A way of making Europe” from the European Union; and by the agreement  
292 OA17/008 with Instituto Tecnológico y de Energías Renovables (ITER) to strengthen

293 scientific and technological education, training, research, development and innovation  
294 in Genomics, Personalized Medicine and Biotechnology.

295

## 296 **Authors' contributions**

297 T.H.B.: performed experiments, data analysis, interpretation, and manuscript drafting;  
298 H.R.P., A.D.U., R.G.M., J.L.S.: performed experiments and data analysis; F.L.D.:  
299 performed experiments, study conception, interpretation, revising the manuscript  
300 critically; C.F.: Study conception and design, data analysis, interpretation, revising the  
301 manuscript critically and conceived the project.

302 The authors declare that there are no competing interests.

303

## 304 **Data availability**

305 All raw reads from MinION and Illumina are available from the SRA database  
306 (accession number(s) SRP141332, SRP141319).

307

308

## 309 **References**

- 310 1. Zubair S, de Villiers EP, Younan M, Andersson G, Tettelin H, Riley DR, et al. Genome  
311 Sequences of Two Pathogenic *Streptococcus agalactiae* Isolates from the One-Humped  
312 Camel *Camelus dromedarius*. *Genome Announc.* 2013;1. doi:10.1128/genomeA.00515-13
- 313 2. Zubair S, de Villiers EP, Fuxelius HH, Andersson G, Johansson K-E, Bishop RP, et al.  
314 Genome Sequence of *Streptococcus agalactiae* Strain 09mas018883, Isolated from a  
315 Swedish Cow. *Genome Announc.* 2013;1. doi:10.1128/genomeA.00456-13
- 316 3. Kropp KA, Lucid A, Carroll J, Belgrudov V, Walsh P, Kelly B, et al. Draft Genome  
317 Sequence of a *Streptococcus agalactiae* Strain Isolated from a Preterm Neonate Blood  
318 Sepsis Patient at the Royal Infirmary, Edinburgh, Scotland. *Genome Announc.* 2014;2.  
319 doi:10.1128/genomeA.00875-14
- 320 4. Areechon N, Kannika K, Hirono I, Kondo H, Unajak S. Draft Genome Sequences of  
321 *Streptococcus agalactiae* Serotype Ia and III Isolates from Tilapia Farms in Thailand.  
322 *Genome Announc.* 2016;4. doi:10.1128/genomeA.00122-16
- 323 5. Barony GM, Tavares GC, Pereira FL, Carvalho AF, Dorella FA, Leal CAG, et al. Large-  
324 scale genomic analyses reveal the population structure and evolutionary trends of  
325 *Streptococcus agalactiae* strains in Brazilian fish farms. *Sci Rep.* 2017;7: 13538.
- 326 6. Bodi Winn C, Dzink-Fox J, Feng Y, Shen Z, Bakthavatchalu V, Fox JG. Whole-Genome  
327 Sequences and Classification of *Streptococcus agalactiae* Strains Isolated from Laboratory-  
328 Reared Long-Evans Rats (*Rattus norvegicus*). *Genome Announc.* 2017;5.  
329 doi:10.1128/genomeA.01435-16
- 330 7. Otaguiri ES, Morguette AEB, Tavares ER, dos Santos PMC, Morey AT, Cardoso JD, et al.  
331 Commensal *Streptococcus agalactiae* isolated from patients seen at University Hospital of  
332 Londrina, Paraná, Brazil: capsular types, genotyping, antimicrobial susceptibility and

- 333 virulence determinants. *BMC Microbiol.* 2013;13: 297.
- 334 8. Flécharde M, Gilot P. Physiological impact of transposable elements encoding DDE  
335 transposases in the environmental adaptation of *Streptococcus agalactiae*. *Microbiology.*  
336 2014;160: 1298–1315.
- 337 9. Morici E, Simoni S, Brenciani A, Giovanetti E, Varaldo PE, Mingoia M. A new mosaic  
338 integrative and conjugative element from *Streptococcus agalactiae* carrying resistance  
339 genes for chloramphenicol (catQ) and macrolides [mef(I) and erm(TR)]. *J Antimicrob*  
340 *Chemother.* Oxford University Press; 2017;72: 64–67.
- 341 10. van der Mee-Marquet N, Diene SM, Barbera L, Courtier-Martinez L, Lafont L, Ouachée  
342 A, et al. Analysis of the prophages carried by human infecting isolates provides new  
343 insight into the evolution of Group B *Streptococcus* species. *Clin Microbiol Infect.* 2017;  
344 doi:10.1016/j.cmi.2017.08.024
- 345 11. Glaser P, Rusniok C, Buchrieser C, Chevalier F, Frangeul L, Msadek T, et al. Genome  
346 sequence of *Streptococcus agalactiae*, a pathogen causing invasive neonatal disease. *Mol*  
347 *Microbiol.* 2002;45: 1499–1513.
- 348 12. Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome  
349 analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the  
350 microbial “pan-genome.” *Proc Natl Acad Sci U S A.* 2005;102: 13950–13955.
- 351 13. Leggett RM, Clark MD. A world of opportunities with nanopore sequencing. *J Exp Bot.*  
352 2017; doi:10.1093/jxb/erx289
- 353 14. Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with  
354 multiplex MinION sequencing [Internet]. *bioRxiv.* 2017. p. 160614. doi:10.1101/160614
- 355 15. Jain M, Fiddes IT, Miga KH, Olsen HE, Paten B, Akeson M. Improved data analysis for  
356 the MinION nanopore sequencer. *Nat Methods.* 2015;12: 351–356.

- 357 16. de Lannoy C, de Ridder D, Risse J. The long reads ahead: de novo genome assembly using  
358 the MinION. *F1000Res*. 2017;6: 1083.
- 359 17. Norris AL, Workman RE, Fan Y, Eshleman JR, Timp W. Nanopore sequencing detects  
360 structural variants in cancer. *Cancer Biol Ther*. 2016;17: 246–253.
- 361 18. Stancu MC, van Roosmalen MJ, Renkens I, Nieboer M, Middelkamp S, de Ligt J, et al.  
362 Mapping And Phasing Of Structural Variation In Patient Genomes Using Nanopore  
363 Sequencing [Internet]. *bioRxiv*. 2017. p. 129379. doi:10.1101/129379
- 364 19. Jain M, Olsen HE, Paten B, Akeson M. The Oxford Nanopore MinION: delivery of  
365 nanopore sequencing to the genomics community. *Genome Biol*. 2016;17: 239.
- 366 20. Rand AC, Jain M, Eizenga JM, Musselman-Brown A, Olsen HE, Akeson M, et al.  
367 Mapping DNA methylation with high-throughput nanopore sequencing. *Nat Methods*.  
368 2017;14: 411–413.
- 369 21. Byrne A, Beaudin AE, Olsen HE, Jain M, Cole C, Palmer T, et al. Nanopore Long-Read  
370 RNAseq Reveals Widespread Transcriptional Variation Among the Surface Receptors of  
371 Individual B cells [Internet]. *bioRxiv*. 2017. p. 126847. doi:10.1101/126847
- 372 22. Castro-Wallace SL, Chiu CY, John KK, Stahl SE, Rubins KH, McIntyre ABR, et al.  
373 Nanopore DNA Sequencing and Genome Assembly on the International Space Station  
374 [Internet]. *bioRxiv*. 2016. p. 077651. doi:10.1101/077651
- 375 23. Márquez S, Carrera J, Pullan ST, Lewandowski K, Paz V, Loman N, et al. First Complete  
376 Genome Sequences of Zika Virus Isolated from Febrile Patient Sera in Ecuador. *Genome*  
377 *Announc*. 2017;5. doi:10.1128/genomeA.01673-16
- 378 24. Chandler J, Camberis M, Bouchery T, Blaxter M, Le Gros G, Eccles DA. Annotated  
379 mitochondrial genome with Nanopore R9 signal for *Nippostrongylus brasiliensis*.  
380 *F1000Res*. 2017;6: 56.



- 381 25. Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using  
382 only nanopore sequencing data. *Nat Methods*. 2015;12: 733–735.
- 383 26. Sovic I, Krizanovic K, Skala K, Sikic M. Evaluation of hybrid and non-hybrid methods for  
384 de novo assembly of nanopore reads [Internet]. *bioRxiv*. 2015. p. 030437.  
385 doi:10.1101/030437
- 386 27. Salazar AN, Gorter de Vries AR, van den Broek M, Wijsman M, de la Torre Cortés P,  
387 Brickwedde A, et al. Nanopore sequencing enables near-complete de novo assembly of  
388 *Saccharomyces cerevisiae* reference strain CEN.PK113-7D. *FEMS Yeast Res*. 2017;17.  
389 doi:10.1093/femsyr/fox074
- 390 28. Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, et al. Real-time,  
391 portable genome sequencing for Ebola surveillance. *Nature*. 2016;530: 228–232.
- 392 29. Mojarro A, Hachey J, Ruvkun G, Zuber MT, Carr CE. CarrierSeq: a sequence analysis  
393 workflow for low-input nanopore sequencing [Internet]. *bioRxiv*. 2017. p. 175281.  
394 doi:10.1101/175281
- 395 30. Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz MC, McCombie WR.  
396 Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a  
397 eukaryotic genome. *Genome Res*. 2015;25: 1750–1756.
- 398 31. Karlsson E, Lärkeryd A, Sjödin A, Forsman M, Stenberg P. Scaffolding of a bacterial  
399 genome using MinION nanopore sequencing. *Sci Rep*. 2015;5: 11996.
- 400 32. Batovska J, Lynch SE, Rodoni BC, Sawbridge TI, Cogan NO. Metagenomic arbovirus  
401 detection using MinION nanopore sequencing. *J Virol Methods*. 2017;249: 79–84.
- 402 33. Risse J, Thomson M, Patrick S, Blakely G, Koutsovoulos G, Blaxter M, et al. A single  
403 chromosome assembly of *Bacteroides fragilis* strain BE1 from Illumina and MinION  
404 nanopore sequencing data. *Gigascience*. 2015;4: 60.

- 405 34. Fournier T, Gounot J-S, Freel K, Cruaud C, Lemainque A, Aury J-M, et al. High-Quality  
406 de Novo Genome Assembly of the *Dekkera bruxellensis* Yeast Using Nanopore MinION  
407 Sequencing. *G3* . 2017;7: 3243–3250.
- 408 35. Jansen HJ, Liem M, Jong-Raadsen SA, Dufour S, Weltzien F-A, Swinkels W, et al. Rapid  
409 de novo assembly of the European eel genome from nanopore sequencing reads. *Sci Rep*.  
410 2017;7: 7213.
- 411 36. Lorenzo-Diaz F, Fernández-Lopez C, Douarre P-E, Baez-Ortega A, Flores C, Glaser P, et  
412 al. Streptococcal group B integrative and mobilizable element IMESag-rpsI encodes a  
413 functional relaxase involved in its transfer. *Open Biol*. 2016;6. doi:10.1098/rsob.160084
- 414 37. Da Cunha V, Davies MR, Douarre P-E, Rosinski-Chupin I, Margarit I, Spinali S, et al.  
415 *Streptococcus agalactiae* clones infecting humans were selected and fixed through the  
416 extensive use of tetracycline. *Nat Commun*. 2014;5: 4544.
- 417 38. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.  
418 *J Mol Biol*. 1990;215: 403–410.
- 419 39. Stewart RD, Watson M. poRe GUIs for parallel and real-time processing of MinION  
420 sequence data. *Bioinformatics*. 2017;33: 2207–2208.
- 421 40. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome  
422 assemblies from short and long sequencing reads. *PLoS Comput Biol*. 2017;13: e1005595.
- 423 41. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome  
424 assemblies. *Bioinformatics*. 2013;29: 1072–1075.
- 425 42. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo  
426 genome assemblies. *Bioinformatics*. 2015;31: 3350–3352.
- 427 43. Baez-Ortega A, Lorenzo-Diaz F, Hernandez M, Gonzalez-Vila CI, Roda-Garcia JL,  
428 Colebrook M, et al. IonGAP: integrative bacterial genome analysis for Ion Torrent

- 429 sequence data. *Bioinformatics*. 2015;31: 2870–2873.
- 430 44. Darling AE, Tritt A, Eisen JA, Facciotti MT. Mauve assembly metrics. *Bioinformatics*.  
431 2011;27: 2756–2757.
- 432 45. George S, Pankhurst L, Hubbard A, Votintseva A, Stoesser N, Sheppard AE, et al.  
433 Resolving plasmid structures in Enterobacteriaceae using the MinION nanopore sequencer:  
434 assessment of MinION and MinION/Illumina hybrid data assembly approaches. *Microbial*  
435 *Genomics*. 2017;3. doi:10.1099/mgen.0.000118
- 436 46. Greninger AL, Naccache SN, Federman S, Yu G, Mbala P, Bres V, et al. Rapid  
437 metagenomic identification of viral pathogens in clinical samples by real-time nanopore  
438 sequencing analysis. *Genome Med*. 2015;7: 99.
- 439 47. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and  
440 accurate long-read assembly via adaptivek-mer weighting and repeat separation. *Genome*  
441 *Res*. 2017;27: 722–736.
- 442 48. Quick J, Quinlan AR, Loman NJ. A reference bacterial genome dataset generated on the  
443 MinION™ portable single-molecule nanopore sequencer. *Gigascience*. 2014;3: 22.
- 444 49. Jain M, Koren S, Quick J, Rand AC, Sasani TA, Tyson JR, et al. Nanopore sequencing and  
445 assembly of a human genome with ultra-long reads [Internet]. *bioRxiv*. 2017. p. 128835.  
446 doi:10.1101/128835
- 447 50. Judge K, Hunt M, Reuter S, Tracey A, Quail MA, Parkhill J, et al. Comparison of bacterial  
448 genome assembly software for MinION data and their applicability to medical  
449 microbiology. *Microb Genom*. 2016;2: e000085.
- 450 51. Benítez-Páez A, Portune KJ, Sanz Y. Species-level resolution of 16S rRNA gene  
451 amplicons sequenced through the MinION™ portable nanopore sequencer. *Gigascience*.  
452 2016;5: 4.

453 52. Shin J, Lee S, Go M-J, Lee SY, Kim SC, Lee C-H, et al. Analysis of the mouse gut  
454 microbiome using full-length 16S rRNA amplicon sequencing. *Sci Rep.* 2016;6: 29681.

455 53. Zaaier S, Columbia University Ubiquitous Genomics 2015 class, Erlich Y. Using mobile  
456 sequencers in an academic classroom. *Elife.* 2016;5. doi:10.7554/eLife.14258.

457

458

459

460

461

462 **Fig 1.** Schematic representation of the *S. agalactiae* HRC hybrid assembly. The  
463 flanking sequences of the two characteristics elements of this isolate are indicated at the  
464 bottom: 1) A conjugative element attached to the Tn916 transposon located in the  
465 largest contig (1,018,988 bp); 2) A prophage located in the third largest contig (224,154  
466 bp). Most copies of the rRNA (16S, 23S, 5S) and tRNA genes were positioned within  
467 the largest unresolved assembly structure.

468

469

470 **Fig 2.** Distribution of MinION reads by quality (Q) scores and length (bp). Aligned  
471 reads are represented in grey, whereas the 68 unaligned reads are overlaid in black. For  
472 simplicity, the plot was capped to represent reads <60 kb length and Q scores <15.

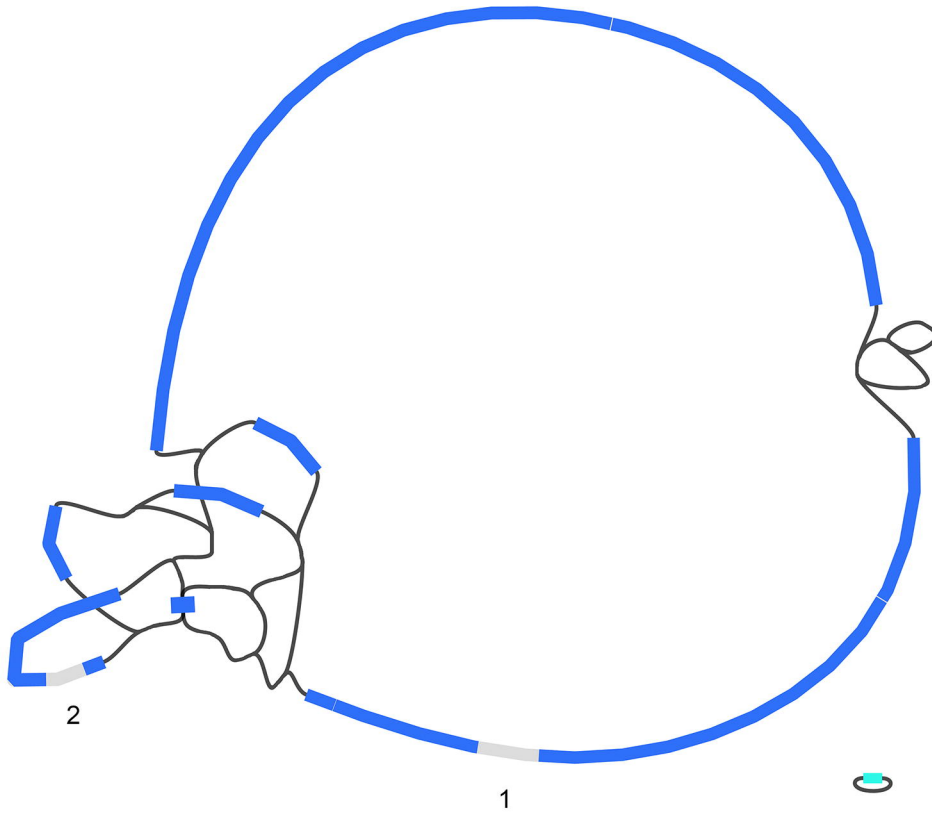
473

**Table 1.** Summary of the hybrid and Illumina-only assemblies of *S. agalactiae* HRC genome.

Parameters	Illumina	Hybrid
Total length assembled (bp)	2,144,732	2,159,060
Contigs (>=500 bp)	23	8
Largest contig (bp)	525,288	1,018,988
N50	147,211	714,293
GC content (%)	35.3	35.4
Genome fraction (%)*	97.9	98.5
Total aligned length (bp)*	2,049,610	2,062,825
Largest alignment (bp)*	524,861	1,018,010
Mismatches (#)*	151	193
Indels (#)*	33	77
Genes**	2135	2153
Coding DNA sequences**	2104	2117
tRNAs**	31	36

\*Against the reference sequence: NZ\_CP010867.1.

\*\*For comparison, IonGAP predicted 2,195 genes, 2,115 coding DNA sequences, and 80 tRNAs in the reference.



1 ...CAATGACACATCATGTCGAGACGGTAGCACTTTTGTCCAA...TGGACAATGATGAACAAGTTGAGTGTGTTGCTCTGCTTGT...

2 ...GCATTACACTTTTAGAAATCAAGGATAGTAAATTTCTTT...TAATTTTAGTAACCTTGATTGTGGTTAGTGTAGAGCCTT...

