

1 **Evaluation of Nanopore sequencing technology to differentiate *Salmonella* serotypes and**

2 **serotype variants with the same or closely related antigenic formulae**

3 Feng Xu^{1,*}, Chongtao Ge^{1,*}, Shaoting Li², Silin Tang¹, Xingwen Wu¹, Hao Luo¹, Xiangyu

4 Deng², Guangtao Zhang¹, Abigail Stevenson¹, Robert C. Baker¹

5 ¹ Mars Global Food Safety Center, Beijing, 101407, China

6 ² Center for Food Safety, University of Georgia, Griffin, GA, 30223, USA

7 *Corresponding authors: feng.y.xu@effem.com; chongtao.ge@effem.com

8 **Highlights**

- 9 • *Salmonella* serotypes or serotype variants with the same antigenic formula were
10 differentiated by SNP typing.
- 11 • Nanopore sequencing followed by phage prediction identified the *Salmonella*
12 serotype variants caused by phage conversion.
- 13 • The latest ONT technology is capable of high fidelity SNP typing of *Salmonella*.

14 **ABSTRACT:**

15 Our previous study demonstrated that whole genome sequencing (WGS) data generated by
16 Oxford Nanopore Technologies (ONT) can be used for rapid and accurate prediction of
17 *Salmonella* serotypes. However, one limitation is that established methods for WGS-based
18 serotype prediction cannot differentiate certain serotypes and serotype variants with the same
19 or closely related antigenic formulae. This study aimed to evaluate Nanopore sequencing and
20 corresponding data analysis for differentiation of these serotypes and serotype variants, thus
21 overcoming this limitation. Five workflows that combined different flow cells, library
22 construction methods and basecaller models were evaluated and compared. The workflow that
23 consisted of the R9 flow cell, rapid sequencing library construction kit and guppy basecaller
24 with base modified model performed best for Single Nucleotide Polymorphism (SNP)
25 analysis. With this workflow, as high as 99.98% matched the identity of the assembled
26 genomes and only less than five high quality SNPs (hqSNPs) between ONT and Illumina
27 sequencing data were achieved. SNP typing allowed differentiation of *Choleraesuis sensu*
28 *stricto*, *Choleraesuis* var. *Kunzendorf*, *Choleraesuis* var. *Decatur*, *Paratyphi C*, and *Typhisuis*
29 that share the same antigenic formula 6,7:c:1,5. Prophage prediction further distinguished
30 Orion var. 15⁺ and Orion var. 15⁺, 34⁺. Our study improves the readiness of ONT as a
31 *Salmonella* subtyping and source tracking tool for food industry applications.

32 **Keywords:** *Salmonella*, Oxford Nanopore sequencing, phylogenetic analysis, SNP, variants

33 **1. Introduction**

34 *Salmonella* remains an important concern in the food industry since this foodborne pathogen
35 continues to cause global public health and economic impact. Thus, it is imperative to
36 reinforce *Salmonella* control measures in the food industry (GMA, 2009). Rapid methods for
37 subtyping strains beyond the species level to serotype are essential to facilitate contamination
38 incident investigations (Olaimat and Holley, 2012; Shi et al., 2015). There are more than
39 2,600 *Salmonella* serotypes currently described in the White-Kauffmann-Le Minor scheme
40 (Grimont and Weill, 2007; Dieckmann and Malorny, 2011). Whole Genome Sequencing
41 (WGS) based tools have now emerged as an alternative method for *Salmonella* serotyping
42 (Zhang et al., 2015; Yoshida et al., 2016; Yachison et al., 2017; Ibrahim and Morin, 2018;
43 Diep et al., 2019; Tang et al., 2019; Zhang et al., 2019; Uelze et al., 2020) and have been
44 suggested to be the new gold standard (Banerji et al., 2020).

45 Oxford Nanopore Technologies (ONT) sequencing generates long-read sequences with
46 rapid turnaround and has the potential for rapid identification of food-borne pathogens during
47 routine monitoring and incident investigations in the food industry. Challenges in the
48 basecalling accuracy of nanopore sequencing have been identified compared with short-read
49 sequencing technologies (Amarasinghe et al., 2020). Errors caused by signal interpretation
50 during basecalling are dominated by insertions and deletions (indels), especially in
51 homopolymeric regions (Gargis et al, 2019). An average of 95% accuracy of raw read (Jain et
52 al., 2017, 2018) and 1.06 errors per Kbp assembly (Lang et al., 2020) were recently reported.

53 *In silico* serotype prediction using ONT sequences was systematically evaluated in our
54 previous study, and all the tested *Salmonella* isolates representing 34 serotypes were correctly
55 predicted at the level of antigenic profile (Xu et al., 2020). Comparison with the conventional

56 serotyping method using antisera combined with biochemical tests, identified a limitation that
57 the established WGS methods, with either Illumina or ONT data, were unable to differentiate
58 serotypes that share the same or related antigenic formulae.

59 Serotypes with antigenic formula 6,7:c:1,5 including Paratyphi C (Vi + or Vi-), Typhisuis,
60 Choleraesuis *sensu stricto*, Choleraesuis var. Kunzendorf and Choleraesuis var. Decatur,
61 require additional biochemical tests to differentiate between them in conventional serotyping
62 (Grimont and Weill, 2007). Choleraesuis is a serotype adapted to swine. This serotype has
63 caused serious disease outbreaks in swine, including two in Denmark in 1999 - 2000 and 2012
64 - 2013 (Leekitcharoenphon et al., 2019). It also has a propensity to cause extraintestinal
65 infections in humans (Chiu et al., 2004; Sirichote et al., 2010). Choleraesuis var. Kunzendorf
66 is responsible for the majority of outbreaks among swine (Leekitcharoenphon et al., 2019).
67 Serotype Orion has two closely related variants, Orion var. 15⁺ and Orion var. 15⁺, 34⁺. Phage
68 conversion by ϕ_{15} or ϕ_{34} is needed to differentiate these two variants (Grimont and Weill,
69 2007). Orion has been isolated from a variety of hosts worldwide, including humans (Cabrera
70 et al., 2006; Trafny et al., 2006), cattle (Alam et al., 2009), ducks (Rampersad et al., 2008),
71 dogs (Gupta et al., 2016) and birds (Münch et al., 2012). Orion var.15⁺ and Orion var. 15⁺,
72 34⁺ are more likely to be associated with poultry (Corry et al., 2002; McWhorter et al., 2015).
73 Additional subtyping analysis is needed to differentiate these variants.

74 When the genetic determinant for a particular biotype or variant is known, it makes an
75 effective target for further differentiation. For example, SeqSero2 uses a Single Nucleotide
76 Polymorphism (SNP) that inactivates tartrate fermentation to differentiate the *S. Paratyphi B*
77 pathogen and a 7-bp deletion to identify O5- variants of *S. Typhimurium* (Zhang et al., 2019).

78 When such genetic determinants are not available, phylogenetic analysis, such as SNP typing
79 and core/whole genome Multilocus Sequence Typing (MLST), can be used to differentiate
80 biotypes and serotype. Nair et al. grouped the 6,7:c:1,5 biotypes of *S. enterica* subspecies
81 *enterica* including Paratyphi C, Typhisuis, Choleraesuis *sensu stricto*, Choleraesuis var.
82 Kunzendorf into four phyloclusters using the core SNP phylogeny (Nair et al., 2020). The
83 authors also found that Choleraesuis var. Decatur genomes showed more SNP variation
84 between them than was found across all of the remaining biotypes (Nair et al., 2020). With
85 continuous improvement of the sequencing and basecalling accuracy of ONT, application of
86 ONT sequences in SNP analysis has become more applicable (Greig et al., 2019; Taylor et al.,
87 2019).

88 The goal of this study was to evaluate the combined use of latest improved ONT
89 sequencing and modified subtyping analyses with primary serotype prediction in
90 differentiating representative *Salmonella* serotypes or serotype variants that shared the same
91 or related antigenic formulae.

92 **2. Materials and methods**

93 *2.1. Bacterial strains*

94 Two strains, *S. enterica* Choleraesuis var. Kunzendorf and *S. enterica* Orion var. 15⁺, 34⁺,
95 were obtained from Dr. Martin Weidmann's laboratory (Cornell University, USA). Details of
96 these strains can be found at www.foodmicrobetracker.com under the isolate IDs FSL
97 R9-0095 and FSL R8-3858.

98 *2.2. Genomic DNA extraction and quality control*

99 Strains were incubated for 20 hrs on Trypticase Soy Agar (TSA) at 37°C. Genomic DNA was

100 extracted using a QIAamp DNA mini kit (Qiagen, Hilden, Germany), following the
101 instructions provided by the manufacture. The double-stranded DNA (dsDNA) was quantified
102 using a Qubit 3.0 fluorimeter (Life Technologies, Paisley, UK). The quality of DNA was
103 assessed using a NanoDrop 1000 instrument (Thermo Fisher Scientific, Delaware, USA).

104 *2.3. Illumina sequencing and genome assembly*

105 Illumina sequences were performed at the Beijing Genomics Institute (BGI, Shenzhen, China).
106 The library for sequencing was constructed by BGI and sequenced using a HiSeq X Ten
107 platform (Illumina, San Diego, CA, USA) to 2×150 cycles. Genomes were assembled by the
108 SOAPdenovo v1.05 (URL: <http://soap.genomics.org.cn>).

109 *2.4. Oxford Nanopore sequencing and genome assembly*

110 Five workflows with combinations of different flow cells, library construction methods and
111 basecaller models were used in this study for comparison purposes (Supplementary Table 1).
112 The DNA library was prepared with library construction kits following the manufacturer's
113 instructions as indicated in each workflow. For flow cell, compared with commercial
114 available R9 (version 9.4.1), R10 (version 10.0, early access) is a new design of nanopore,
115 with a longer barrel and dual reader head. Libraries were sequenced with qualified R9 and
116 R10 flow cells (active pores number ≥ 800) on a GridION sequencer (Oxford Nanopore
117 Technologies, Oxford, UK). Real-time basecalling was performed using Guppy version 3.2.6
118 with the corresponding basecalling model as indicated in each workflow, which was
119 integrated in the MinKNOW software v3.5.4. Fastq data were obtained for further analysis.
120 Adaptor sequences were trimmed with Porechop v0.2.3 (URL:

121 <https://github.com/rrwick/Porechop>), and the quality of trimmed data was assessed using
122 NanoStat v1.1.2 (De Coster et al., 2018). Filtlong v0.2.0 (URL:
123 <https://github.com/rrwick/Filtlong>) was applied to remove reads shorter than 1,000 bp.
124 Genomes were *de novo* assembled by Wtdbg2 v2.4 (Ruan and Li, 2019) with default
125 parameters. The assembled genomes were corrected using Racon v1.3.3 (Vaser et al. 2017).
126 Consensus was obtained from one round of Medaka version 0.8.0
127 (<https://github.com/nanoporetech/medaka>). QUAST v5.0.2 was applied to assess the
128 assembled contigs (Gurevich et al., 2013). DNAdiff v1.3 (MUMmer version, Kurtz et al.,
129 2004) was used to evaluate the base level comparison between ONT and Illumina assemblies.

130 2.5. Phylogenetic analysis

131 2.5.1. kSNP tree

132 A k-mer based approach, kSNP3 (Gardner et al., 2015) was used to assess the clustering of
133 isolates sharing the same antigenic formula with FSL R9-0095 and isolates sharing closely
134 related antigenic formulae with FSL R8-3858. Kchooser was used to determine an optimum
135 k-mer size (Carroll et al., 2017). A parsimony tree was generated by kSNP3 using the set of
136 core SNPs identified. Clade robustness was assessed using a bootstrap analysis with 1,000
137 replicates (Felsenstein, 1985). Twenty-six draft genome sequences were downloaded from
138 GenBank and further confirmed by SeqSero2 (Zhang et al., 2019) and SISTR (Yoshida et al.,
139 2016) for the phylogenetic analysis of strain FSL R9-0095 (Supplementary Table 2). These
140 draft genome sequences included six strains of Typhisuis, eight strains of Paratyphi C, two
141 strains of Choleraesuis *sensu stricto*, three strains of Choleraesuis var. Decatur, and seven

142 strains of *Choeraesuis* var. *Kunzendorf*. The Illumina sequences of strain FSL R9-0095 and
143 the corresponding ONT sequences generated using all the sequencing workflows were
144 analyzed through kSNP analysis. Twenty-seven draft genomes were downloaded from
145 GenBank and further confirmed by SeqSero2 and SISTR for the phylogenetic analysis of
146 strain FSL R8-3858 (Supplementary Table 2). These draft genomes included 16 strains of
147 Orion, seven strains of Orion var. 15⁺, and four strains of Orion var. 15⁺, 34⁺. The Illumina
148 sequences of strain FSL R8-3858 and the corresponding ONT sequences generated using all
149 the sequencing workflows were analyzed through kSNP analysis.

150 2.5.2. *CFSAN high quality SNP (hqSNP) pipeline*

151 The hqSNP pipeline developed by the Center for Food Safety and Applied Nutrition (CFSAN
152 SNP Pipeline v.1.0.0/FDA) was used for SNP calling of *Salmonella* strains with default
153 quality filters (Davis et al., 2015). Read mapping was performed using Burrows-Wheeler
154 Aligner (v0.7.17-r1188) (BWA-MEM) with the -x ont2d option (Li and Durbin, 2010; Hyeon
155 et al., 2018) for the ONT sequences. An SNP matrix and alignment concatenated SNP were
156 produced with customized Python scripts.

157 2.5.3. *Phylogenetic tree generated from CFSAN hqSNP analysis*

158 Maximum likelihood (ML) phylogenetic trees were built using PhyML v3.3 based on the
159 hqSNP matrices and visualized with Figtree v1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>).
160 All the phylogenetic trees shown in this paper were midpoint rooted.

161 2.6. *Phage detection*

162 All the Illumina data of Orion and its varirants, as well as the assembled genomes of strain

163 FSL R8-3858 sequenced using either Illumina or ONT, were submitted to the Phage Search
164 Tool - Enhanced Release (PHASTER API website, <http://phaster.ca>). Results were extracted
165 from the files returned from the server (Zhou et al., 2011; Arndt et al., 2016; Worley et al.,
166 2018). Sequences classified by PHASTER as questionable (score between 70 and 90), intact
167 (score > 90) or incomplete (score < 70) were all recorded.

168 **3. Results and discussion**

169 *3.1. Overview of ONT raw reads*

170 The same ONT fast5 data were used for basecalling by the high accuracy (HAC) model and
171 base modified model (workflow 1 vs workflow 2) (Supplementaryable 1). The data with
172 quality score ≥ 12 were considered high quality data. The mean read length, mean read
173 quality and percentage of high-quality data (high quality data/total bases) were similar
174 between data generated by these two basecalling models (Table 1). The read length of
175 sequences generated from workflow 3 was close to 4 Kbp (Table 1). One strand of dsDNA
176 prepared by the 1D2 library preparation kit was sequenced followed by its complementary
177 strand for workflow 4. This workflow generated sequences with the longest mean read length
178 (Table 1). However, its mean read quality of the sequences was the lowest (Table 1).

179 Two equal aliquots of the same batch of extracted genomic DNA were sequenced with the
180 R9 and R10 flow cells for workflows 1 and 5, respectively. The yield of the R9 flow cell was
181 higher than that of the R10 flow cell (Table 1). The mean read length of sequences from the
182 R9 flow cell were shorter than that from the R10 flow cell (Table 1). The mean read quality of
183 sequences from the R9 flow cell were higher than that from the R10 flow cell (Table 1). The
184 percentage of high quality data of sequences from the R9 flow cell were higher than that from

185 the R10 flow cell (Table 1).

186 According to the previous study sequences produced in two hrs from ONT were sufficient
187 for *Salmonella* serotype prediction (Xu et al., 2020). Therefore, the raw ONT sequences from
188 the first two hrs of sequencing time were extracted to investigate whether they were sufficient
189 for additional subtyping analysis. About 100× genome coverage was achieved after two hrs
190 sequencing, except with workflow 5 which used the R10 flow cell for FSL R9-0095
191 sequencing (58×, Table 1). This might have resulted from a dramatic decrease of active pores
192 on the R10 from 1,223 to 155 after the library was loaded, although it increased after 1.5 hrs.
193 Overall, comparisons of the sequencing data yield, mean read length, mean quality score and
194 high-quality data percentage from different workflow using two hrs sequencing data were in
195 line with those using the data of full runs (24, 48, 72 hrs, Table 1).

196 *3.2. Overview of ONT assembled contigs*

197 The ONT sequencing data were assembled using Wtdbg2, and corrected by 2 rounds of
198 Racon. Consensus were then obtained by one round of Medaka (Supplementary Figure 1).
199 The overview of the assembled contigs through the reads obtained on GridION are
200 summarized below (Table 2). The contigs generated from Illumina sequencing were used as a
201 benchmark for comparisons. In general, higher matching identity (> 99.90%) and fewer SNPs
202 by DNAdiff (< 1,000) were obtained from workflows 2, 3 and 5 compared with workflows 1
203 or 4 (< 99.90% of matching identity and >1,000 of SNPs by DNAdiff). DNA libraries
204 prepared with the rapid kit (workflows 1, 2 and 5) and 1D2 kit (workflow 4) generated fewer
205 contigs (< 5) and higher N50 length (> 4.6 Mbp) in general, except for the 48 hrs data
206 obtained with workflow 5 for FSL R8-3858 sequencing (N50 = 3.51 Mbp). The DNA library

207 prepared using the PCR kit (workflow 3) generated more contigs and shorter N50 length
208 (Table 2). The reason could be fragmentation of genomic DNA followed by a PCR
209 amplification step during the library preparation. Taylor et al. (2019) reported that four hrs
210 ONT sequences generated full-length genomes with an average identity of 99.87% for *S.*
211 *Bareilly* and 99.89% for *Escherichia coli* in comparison to the respective Illumina references.
212 With the accuracy improvement of ONT sequencing, especially for the basecaller and flow
213 cell, the highest matching identity was improved to 99.99% as demonstrated using sequencing
214 workflow 3 with the PCR kit in this study (Table 2).

215 Two sequencing workflows were tested to reduce the impact of methylation, as an average
216 of 95% of the discrepant positions between the sequencing technologies (ONT *versus*
217 Illumina) were caused by methylations (Greig et al., 2019). These were workflow 2 (R9 flow
218 cell + rapid kit + base modified model basecaller) and workflow 3 (R9 flow cell + PCR kit +
219 HAC model basecaller). The base modified model basecaller was trained with human and *E.*
220 *coli* reads, which allowed for calling 5-methylcytosine (5mC) and N6-methyladenine (6mA)
221 modified bases. In comparison with workflow 1, the total SNP loci generated from workflow
222 2 decreased from 4492 - 4854 to 29 - 38 for both tested strains (Table 2). Since workflow 3
223 included a PCR amplification step in library preparation, it could exclude methylations from
224 the library (Liu et al., 2019). Compared with workflow 1, the total SNP loci decreased to 42 -
225 245 for both tested strains (Table 2). Workflow 4 generated similar SNP numbers to that of
226 workflow 1, ranging from 4,400 to 6,000, which suggested that sequencing both DNA strands
227 might not improve accuracy (Table 2). Despite the low data quantity and genome coverage,
228 fewer SNPs were obtained by sequencing on R10 flow cells (workflow 5) compared to R9

229 flow cells (workflow 1). R10 flow cell had a longer barrel and dual reader head, enabling
230 improved resolution of homopolymeric regions and improving the consensus accuracy of
231 ONT sequencing data
232 (<https://nanoporetech.com/about-us/news/r103-newest-nanopore-high-accuracy-nanopore-sequencing-now-available-store>). Current results indicated that flow cell R10 flow cell did
233 improve sequencing accuracy.
234

235 *3.3. Differentiation of serotypes and serotype variants that share the same antigenic formula*

236 *3.3.1 SNP analysis*

237 Isolates belonging to each of the three serotypes (Typhisuis, Paratyphi C and Choleraesuis)
238 and two serotype variants (Chloeraesuis var. Decatur and Chloeraesuis var. Kunzendorf)
239 sharing the same antigenic formula 6,7:c:1,5 formed a distinct clade in the kSNP phylogeny
240 with 100% bootstrap support for each clade (Figure 1). This suggested that kSNP analysis
241 using assembled genomes could further differentiate serotypes and variants with this antigenic
242 formula. Genomes of strain FSL R9-0095 sequenced with either the Illumina or ONT
243 platforms fell into the Choleraesuis var. Kunzenforf clade (Figure 1), which indicated that
244 ONT sequencing was equivalent to Illumina for kSNP analysis in identifying these serotypes
245 and serotype variants. Since the SNP numbers generated through sequencing workflows 1 and
246 4 were higher than those through the other three workflows, the branches of these four
247 corresponding genomes were longer than the other ONT genomes with 100% of bootstrap
248 support (Figure 1). A shorter branch was observed in the subclade through use of workflow 5
249 (sequencing on the R10 flow cell).

250 Genomes of strain FSL R9-0095 sequenced on the Illumina platform were used as a

251 reference for read mapping in the analysis of CFSAN hqSNP. ONT sequences generated
252 using the five different workflows and from different sequencing time were analyzed using
253 the CFSAN hqSNP pipeline. The numbers of hqSNP are listed below (Table 2). Sequencing
254 workflows 2 and 3 generated 0 hqSNP between Illumina data and corresponding ONT data
255 for strain FSL R9-0095. Workflow 5 generated 1 hqSNP after two or 24 hrs of sequencing of
256 the same strain. Overall, the 5 sequencing workflows yielded fewer than 10 hqSNPs
257 compared with the corresponding Illumina sequences (Table 2).

258 The phylogenetic trees generated by hqSNP using ONT and/or Illumina sequences for
259 strain FSL R9-0095 are presented below (Figure 2a, 2b and 2c). Three phylogenetic trees
260 were constructed. In order to compare ONT and Illumina sequences, the first tree was built
261 using both ONT and Illumina sequences, including Illumina sequences from GenBank,
262 Illumina sequence of the tested isolate, and ONT sequences generated by different workflows
263 and different sequencing time (Figure 2a). In the second tree, one ONT-sequenced isolate was
264 included to replace the corresponding Illumina-sequenced genome to assess hqSNP-based
265 phylogenetic analysis using ONT sequencing (Figure 2b). The third tree was built using
266 Illumina sequences of these isolates only. Strain FSL R9-0095 was correctly placed in the
267 Choleraesuis var. Kunzendorf clade (Figures 2a, 2b and 2c). All the different serotypes or
268 serotype variants under antigenic formula 6,7:c:1,5 were clearly separated into their
269 respective clades, and all the clades were supported by 100% bootstrap. This indicated that
270 the CFSAN hqSNP-based phylogenetic analysis with either ONT or Illumina data was able to
271 differentiate Choleraesuis var. Kunzendorf isolates. In addition, less than 30 mins was needed
272 for library preparation by the rapid library preparation kit, while about three and five hrs were

273 required using the 1D2 and PCR kits, respectively. In summary, workflow 2 (which used the
274 rapid library preparation kit) was the more optimal procedure to deliver the desired
275 phylogenetic analysis in terms of accuracy and turnaround time. The phylogenetic trees
276 constructed from data sets of two sequencing times (two hrs and 24 hrs) using workflow 2
277 showed identical topology (data not shown), suggesting that ONT sequences collected within
278 two hrs were sufficient for CFSAN hqSNP-based phylogenetic analysis to generate correct
279 variant prediction.

280 Longo et al. (2019) showed that WGS followed by Short Read Sequence Typing for
281 Bacterial Pathogens (SRST2) for MLST (Inouye et al., 2014) in *Salmonella* TypeFinder
282 (<https://cge.cbs.dtu.dk/services/SalmonellaTypeFinder/>) could correctly identify all the
283 analyzed *S. Choleraesuis* var. Kunzendorf isolates. However, *Salmonella* TypeFinder cannot
284 currently accept raw ONT reads as inputs and few tools are available to identify *Choleraesuis*
285 var. Kunzendorf using ONT sequences. In this study, phylogenetic analysis through kSNP or
286 CFSAN hqSNP analysis using ONT sequences successfully identified the isolate FSL
287 R9-0095 to the variant level, which suggests the feasibility of applying ONT sequencing for
288 differentiation of the serotypes or serotype variants sharing the same formula 6,7:c:1,5. The
289 phylogenetic result from this study was in accordance with the result from Nair et al. (2020)
290 using Illumina sequences.

291 *3.4. Differentiation of serotype variants that differ by minor antigens*

292 *3.4.1 SNP analysis*

293 Similar results were obtained from both kSNP and CFSAN hqSNP trees generated using
294 either the ONT or the Illumina-sequenced genomes of FSL R8-3858. All the ONT sequences

295 were clustered with their corresponding Illumina sequences (Supplementary Figure 2 and
296 Supplementary Figure 3a). Serotype Orion and its variants Orion var. 15⁺ and Orion var. 15⁺,
297 34⁺ were indistinguishable in both SNP trees (Supplementary Figure 2 and Supplementary
298 Figure 3), indicating that phylogenetic clustering is not able to differentiate the two Orion
299 variants from each other and Orion.

300 3.4.2. Detection of prophage

301 Since the Orion variants are caused by phage conversion and phylogenetic analysis could not
302 differentiate them, *in silico* prophage detection was used for their differentiation. Prophage
303 prediction using Illumina data correctly identified Orion, Orion var. 15⁺, or Orion var. 15⁺,
304 34⁺ in 27 of 42 isolates (Supplementary Table 1). Four sequences were untypable or
305 questionable since \square_{15} and/or \square_{34} could not be identified (Supplementary Table 1). Moreover,
306 11 isolates had questionable prophage prediction that was inconsistent with the original phage
307 typing records (Supplementary Table 1). A 25.6 Kbp region was detected with Illumina data
308 as phage \square_{15} with a score of 20 (marked as incomplete), and a 36.8 Kbp region was detected
309 as phage \square_{34} with a score of 80 (marked as questionable). Sub-optimal assembly and
310 identification of the two phages was likely due to frequent occurrence of sequencing gaps in
311 the prophage regions caused by short Illumina reads.

312 Twenty-seven of the 42 Illumina sequences tested were identified to be in accordance with
313 their original record through PHASTER (supplementary Table 1). By contrast, both
314 prophages for genomes assembled from ONT data were detected with higher scores (100 -
315 150 for \square_{15} and 80 - 90 for \square_{34}), except the genome assembled from 48 hr sequences using
316 workflow 4 (Table 3). These results showed that the much longer ONT reads (compared with

317 Illumina reads) allowed the detection of more complete prophages for definitive
318 differentiation of Orion variants.

319 **5. Conclusion**

320 In this study we demonstrated that long-read nanopore sequencing technology can be used as
321 a subtyping tool to differentiate *Salmonella* serotypes or serotype variants which share the
322 same or related antigenic formulae. The matching identity between ONT and Illumina
323 sequences was improved to 99.98 – 99.99% using the workflow including the bioinformatics
324 pipeline. SNP-based phylogenetic analysis successfully identified serotypes and serotype
325 variants including Choleraesuis, Choleraesuis var. Kunzendorf, Choleraesuis var. Decatur,
326 Paratyphi C, Typhisuis by using ONT sequences. ONT sequencing can also successfully
327 differentiate variants of Orion, Orion var. 15⁺ and Orion var. 15⁺, 34⁺ by addition of
328 PHASTER prediction.

329 **Acknowledgement**

330 The authors would like to thank Dr. Martin Weidmann at Cornell University for providing
331 strains used in this study. The authors would also like to thank Oxford Nanopore
332 Technologies (Thomas Bray, Lin Jin, Xinran Yu, Xiang Chen, Iain MacLaren, Richard
333 Compton, Stephen Rudd, David Dai) for supporting the establishment of ONT capability at
334 the Mars Global Food Safety Center.

335 **References**

336 Alam, M.J., Renter, D.G., Ives, S.E., Thomson, D.U., Sanderson, M.W., Hollis, L.C.,
337 Nagaraja, T.G., 2009. Potential associations between fecal shedding of *Salmonella* in
338 feedlot cattle treated for apparent respiratory disease and subsequent adverse health

- 339 outcomes. *Vet. Res.* 40, 2. <http://dx.doi.org/10.1051/vetres:2008040>.
- 340 Amarasinghe, S.L., Su, S., Dong, X., Zappia, L., Ritchie M.E., Gouil ,Q., 2020. Opportunities
341 and challenges in long-read sequencing data analysis. *Genome Biol.* 21, 30.
342 <https://doi.org/10.1186/s13059-020-1935-5>.
- 343 Arndt, D., Grant, J.R., Marcu, A., Sajed, T., Pon, A., Liang, Y., Wishart, D.S., 2016.
344 PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res.*
345 44, 16 - 21. <https://doi.org/10.1093/nar/gkw387>.
- 346 Banerji, S., Simon, S., Tille, A., Fruth, A., Flieger, A., 2020. Genome-based *Salmonella*
347 serotyping as the new gold standard. *Sci. Rep.* 10 (1), 4333. Doi:
348 10.1038/s41598-020-61254-1.
- 349 Cabrera, R., Ruiz, J., Ramírez, M., Bravo, L., Fernández, A., Aladueña, A., Echeíta, A.,
350 Gascón, J., Alonso, P.L., Vila, J., 2006. Dissemination of *Salmonella enterica* serotype
351 Agona and multidrug-resistant *Salmonella enterica* serotype Typhimurium in Cuba. *Am.*
352 *J. Trop. Med. Hyg.* 74, 1049 - 1053.
- 353 Carroll, L.M., Wiedmann, M., den Bakker, H., Siler, J., Warchocki, S., Kent, D., Lyalina, S.,
354 Davis, M., Sisco, W., Besser, T., Warnick, L.D., Pereira, R.V., 2017. Whole-genome
355 sequencing of drug resistant *Salmonella enterica* isolates from dairy cattle and humans in
356 New York and Washington States reveals source and geographic associations. *Appl.*
357 *Environ. Microbiol.* 83, e00140-17. <https://doi.org/10.1128/AEM.00140-17>.
- 358 Chiu, C.H., Su, L.H., Chu, C., 2004. *Salmonella enterica* serotype Choleraesuis:
359 epidemiology, pathogenesis, clinical disease, and treatment. *Clin. Microbiol. Rev.* 17,
360 311 - 322. Doi: 10.1128/cmr.17.2.311-322.2004.

- 361 Corry, J.E.L., Allen, V.M., Hudson, W.R., Breslin, M.F., Davies, R.H., 2020. Sources of
362 *Salmonella* on broiler carcasses during transportation and processing: modes of
363 contamination and methods of control. *J. Appl. Microbiol.* 92(3): 424 – 432.
- 364 Davis, S., Pettengill, J.B., Luo, Y., Payne, J., Shpuntoff, A., Rand, H., Strain, E., 2015.
365 CFSAN SNP Pipeline: an automated method for constructing SNP matrices from
366 next-generation sequence data. *PeerJ Computer Science* 1, e20.
367 <https://doi.org/10.7717/peerj-cs.20>.
- 368 De Coster, W., D'Hert, S., Schultz, D.T., Cruts, M., Van Broeckhoven, C., 2018. NanoPack:
369 visualizing and processing long-read sequencing data. *Bioinformatics* 34 (15), 2666 -
370 2669. Doi: 10.1093/bioinformatics/bty149.
- 371 Dieckmann, R., Malorny, B., 2011. Rapid screening of epidemiologically important
372 *Salmonella enterica* subsp. *enterica* serovars by whole-cell matrix-assisted laser
373 desorption ionization - time of flight mass spectrometry. *Appl. Environ. Microbiol.*
374 77(12): 4136 - 4146.
- 375 Diep, B., Barretto, C., Portmann, A-C., Fournier, C., Karczmarek, A., Voets, G., Li, S., Deng,
376 X., Klijjn, A., 2019. *Salmonella* serotyping: comparison of the traditional method to a
377 microarray-based method and an in silico platform using whole genome sequencing data.
378 *Front. Microbiol.* 10, 2554. Doi: 10.3389/fmicb.2019.02554.
- 379 Felsenstein, J., 1985. Confidence-delimitation on phylogenies - an approach using the
380 bootstrap. *Evolution.* 39, 783-791.

- 381 Gardner, S.N., Hall, B.G., 2013. When whole-genome alignments just won't work: kSNP v2
382 software for alignment-free SNP discovery and phylogenetics of hundreds of microbial
383 genomes. PLoS One 8, e81760. <https://doi.org/10.1371/journal.pone.0081760>.
- 384 Gargis, A.S., Cherney, B., Conley, A.B., McLaughlin, H.P., Sue, D., 2019. Rapid detection of
385 genetic engineering, structural variation, and antimicrobial resistance markers in bacterial
386 biothreat pathogens by nanopore sequencing. Sci. Rep. 9(1), 13501. Doi:
387 10.1038/s41598-019-49700-1.
- 388 GMA, 2009. The association of food, beverage and consumer products companies (GMA):
389 Control of *Salmonella* in low-moisture foods.
- 390 Greig, D.R., Jenkins, C., Gharbia, S., Dallman, T.J., 2019. Comparison of single-nucleotide
391 variants identified by Illumina and Oxford Nanopore Technologies in the context of a
392 potential outbreak of Shiga toxin-producing *Escherichia coli*. GigaScience 8, 1 - 12.
393 <https://doi.org/10.1093/gigascience/giz104>.
- 394 Grimont, P., Weill, F., 2007. Antigenic formulae of the *Salmonella* serovars, (9th ed.) Paris:
395 WHO Collaborating Centre for Reference and Research on *Salmonella*.
- 396 Gupta, S.K., McMillan, E.A., Jackson, C.R., Desai, P.T., Porwollik, S., McClelland, M., Hiott,
397 L.M., Humayoun, S.B., Frye, J.G., 2016. Draft genome sequence of *Salmonella enterica*
398 subsp. *enterica* serovar Orion strain CRJJGF 00093 (phylum Gammaproteo bacteria).
399 Genome Announc. 4(5), e01063-16. Doi:10.1128/genomeA.01063-16.
- 400 Gurevich, A., Saveliev, V., Vyahhi, N., Tesler, G., 2013. QUASt: quality assessment tool for
401 genome assemblies. Bioinformatics. 29 (8), 1072 - 1075.
402 Doi:10.1093/bioinformatics/btt086.

-
- 403 Hyeon, J-Y., Li, S., Mann, D.A., Zhang, S., Li, Z., Chen, Y., Deng, X., 2018.
404 Quasimetagenomics based and real-time-sequencing-aided detection and subtyping of
405 *Salmonella enterica* from food samples. Appl. Environ. Microbiol. 84, e02340-17.
406 <https://doi.org/10.1128/AEM.02340-17>.
- 407 Ibrahim, G.M., Morin, P.M., 2018. *Salmonella* serotyping using whole genome sequencing.
408 Front. Microbiol. 9, 2993. <https://doi.org/10.3389/fmicb.2018.02993>.
- 409 Inouye, M., Dashnow, H., Raven, L., Schultz, M., Pope, B.J., Tomita, T., Zobel, J., Holt, K.E,
410 2014. SRST2: Rapid genomic surveillance for public health and hospital microbiology
411 labs. Genome Med. 6(11): 1.46.
- 412 Jain, M., Koren, S., Miga, K.H., Quick, J., Rand, A.C., Sasani, T.A., Tyson, J.R., Beggs, A.D.,
413 Dilthey, A.T., Fiddes, I.T., Malla, S., Marriott, H., Nieto, T., O'Grady, J., Olsen, H.E.,
414 Pedersen, B.S., Rhie, A., Richardson, H., Quinlan, A.R., Snutch, T.P., Tee, L., Paten, B.,
415 Phillippy, A.M., Simpson, J.T., Loman, N.J., Loose, M., 2018. Nanopore sequencing and
416 assembly of a human genome with ultra-long reads. Nat. Biotechnol. 36, 338.
417 <https://doi.org/10.1038/nbt.4060>.
- 418 Jain, M., Tyson, J.R., Loose, M., L C Ip, C., Eccles, D., O'Grady, J., Malla, S., Leggett, R.M.,
419 Wallerman, O., Jansen, H., Zalunin, V., Birney, E., Brown, B., Snutch, T., Olsen, H.,
420 2017. MinION analysis and reference consortium: phase 2 data release and analysis of
421 R9.0 chemistry, vol. 6.
- 422 Kurtz, S., Phillippy, A., Delcher, A.L., Smoot, M., Shumway, M., Antonescu, C., Salzberg,
423 S.L., 2004. Versatile and open software for comparing large genomes. Genome Biol. 5,
424 R12.

- 425 Lang, D., Zhang, S., Ren, P., Liang, F., Sun, Z., Meng, G., Tan, Y., Hu, J., Li, X., Lai, Q.,
426 Han, L., Wang, D., Hu, F., Wang, W., Liu, S., 2020. Comparison of the two up-to-date
427 sequencing technologies for genome assembly: HiFi reads of Pacbio Sequel II system
428 and ultralong reads of Oxford Nanopore. bioRxiv. Doi:
429 <https://doi.org/10.1101/2020.02.13.948489>.
- 430 Leekitcharoenphon, P., Sørensen, G., Löfström, C., Battisti, A., Szabo, I., Wasyl, D., Slowey,
431 R., Zhao, S., Brisabois, A., Kornschöber, C., Kärssin, A., Szilárd, J., Cerný, T., Svendsen,
432 C.A., Pedersen, K., Aarestrup, F.M., Hendriksen, R.S., 2019. Cross-border transmission
433 of *Salmonella* Choleraesuis var. Kunzendorf in European pigs and wild boar: infection,
434 genetics, and evolution. *Front. Microbiol.* 10, 179. Doi: 10.3389/fmicb.2019.00179.
- 435 Li, H., Durbin, R., 2010. Fast and accurate long-read alignment with Burrows-Wheeler
436 transform. *Bioinformatics* 26, 589 - 595.
- 437 Liu, Q., Georgieva, D.C., Egli, D. and Wang, K., 2019. NanoMod: a computational tool to
438 detect DNA modifications using Nanopore long-readsequencing data. *BMC Genomics*
439 20, 78. <https://doi.org/10.1186/s12864-018-5372-8>.
- 440 Longo, A., Petrin, S., Mastrorilli, E., Tiengo, A., Lettini, A.A., Barco, L., Ricci, A., Losasso,
441 C., Cibin, V., 2019. Characterizing *Salmonella enterica* serovar Choleraesuis, var.
442 Kunzendorf: a comparative case study. *Front. Vet. Sci.* 6, 316.
443 Doi:10.3389/fvets.2019.00316.
- 444 McWhorter, A.R., Davos, D., Chousalkar, K.K., 2015. Pathogenicity of *Salmonella* strains
445 isolated from egg shells and the layer farm environment in Australia. *Appl. Environ.*
446 *Microbiol.* 81: 405–414. doi:10.1128/AEM.02931-14.

-
- 447 Mottawea, W., Duceppe, M-O., Dupras, A.A., Usongo, V., Jeukens, J., Freschi, L.,
448 Emond-Rheault, J-G., Hamel, J., Kukavica-Ibrulj, I., Boyle, B., Gill, A., Burnett, E.,
449 Franz, E., Arya, G., Weadge, J.T., Gruenheid, S., Wiedmann, M., Huang, H., Daigle, F.,
450 Moineau, S., Bekal, S., Levesque, R.C., Goodridge, L.D. Ogunremi, D., 2018.
451 *Salmonella enterica* prophage sequence profiles reflect genome diversity and can be used
452 for high discrimination subtyping. *Front. Microbiol.* 9, 836. Doi:
453 10.3389/fmicb.2018.00836.
- 454 Münch, S., Braun, P., Wernery, U., Kinne, J., Pees, M., Flieger, A., Tietze, E., Rabsch, W.,
455 2012. Prevalence, serovars, phage types, and antibiotic susceptibilities of *Salmonella*
456 strains isolated from animals in the United Arab Emirates from 1996 to 2009. *Trop.*
457 *Anim. Health Prod.* 44, 1725 - 1738. <http://dx.doi.org/10.1007/s11250-012-0130-4>.
- 458 Nair, S., Fookes, M., Corton, C., Thomson, N.R., Wain, J., Langridge, G.C., 2020. Genetic
459 Markers in *S. Paratyphi C* Reveal Primary Adaptation to Pigs. *Microorganisms*, 8, 657.
460 doi:10.3390/microorganisms8050657.
- 461 Olaimat, A.N., Holley, R.A., 2012. Factors influencing the microbial safety of fresh produce:
462 a review. *Food Microbiol.* 32 (1), 1 - 19. <https://doi.org/10.1016/j.fm.2012.04.016>.
- 463 Rampersad, J., Johnson, J., Brown, G., Samlal, M., Ammons, D., 2008. Comparison of
464 polymerase chain reaction and bacterial culture for *Salmonella* detection in the Muscovy
465 duck in Trinidad and Tobago. *Rev. Panam. Salud. Publica.* 23, 264 - 267.
466 <http://dx.doi.org/10.1590/S1020-49892008000400006>.
- 467 Ruan, J., Li, H., 2019. Fast and accurate long-read assembly with Wtdbg2. 10.1101/530972.

-
- 468 Shi, C., Singh, P., Ranieri, M.L., Wiedmann, M., Moreno Switt, A.I., 2015. Molecular
469 methods for serovar determination of *Salmonella*. Crit. Rev. Microbiol. 41 (3), 309 - 325.
470 <https://doi.org/10.3109/1040841X.2013.837862>.
- 471 Sirichote, P., Hasman, H., Pulsrikarn, C., Schönheyder, H. C., Samulionienė, J.,
472 Pornruangmong, S., Bangtrakulnonth, A., Aarestrup, F.M., Hendriksen, R.S., 2010.
473 Molecular characterization of extended-spectrum cephalosporinase-producing
474 *Salmonella enterica* serovar Choleraesuis isolates from patients in Thailand and
475 Denmark. J. Clin. Microbiol. 48, 883 - 888. Doi: 10.1128/JCM.01792-09.
- 476 Tang, S. Orsi, R.S., Luo, H., Ge, C., Zhang, G., Baker, R.C., Stevenson A., Wiedmann, M.,
477 2019. Assessment and comparison of molecular subtyping and characterization methods
478 for *Salmonella*. Front Microbiol. 10, 1591.
- 479 Taylor, T.L., Volkening, J.D., DeJesus, E., Simmons, M., Dimitrov, K.M., Tillman, G.E.,
480 Suarez, D.L., Afonso, C.L., 2019. Rapid, multiplexed, whole genome and plasmid
481 sequencing of foodborne pathogens using long-read nanopore technology. Sci. Rep. 9 (1),
482 16350. <https://doi.org/10.1038/s41598-019-52424-x>.
- 483 Trafny, E.A., Kozłowska, K., Szpakowska, M., 2006. A novel multiplex PCR assay for the
484 detection of *Salmonella enterica* serovar Enteritidis in human faeces. Lett. Appl.
485 Microbiol. 43, 673 - 679. <http://dx.doi.org/10.1111/j.1472-765X.2006.02007.x>.
- 486 Vaser, R., I. Sovic, N. Nagarajan, and M. Sikic. 2017. Fast and accurate *de novo* genome
487 assembly from long uncorrected reads. Genome Res. 27 (5), 737 - 746. Doi:
488 10.1101/gr.214270.116.

-
- 489 Uelze, L., Borowiak, M, Deneke, C., Szabó, I., Fischer, J., Tausch, S.H., Malorny, B., 2020.
490 Performance and accuracy of four open-source tools for in silico serotyping of
491 *Salmonella* spp. based on whole-genome short-read sequencing data. Appl. Environ.
492 Microbiol. 86, e02265-19. <https://doi.org/10.1128/AEM.02265-19>.
- 493 Worley, J., Meng, J., Allard, M.W., Brown, E.W., Timme, R.E., 2018. *Salmonella enterica*
494 phylogeny based on whole-genome sequencing reveals two new clades and novel
495 patterns of horizontally acquired genetic elements. mBio 9, e02303-18.
496 <https://doi.org/10.1128/mBio.02303-18>.
- 497 Xu, F., Ge, C., Luo, H., Li, S., Wiedmann, M., Deng, X., Zhang, G., Stevenson, A., Baker,
498 R.C., Tang, S., 2020. Evaluation of real-time nanopore sequencing for *Salmonella*
499 serotype prediction. Food Microbiol. 89, 103452. Doi: 10.1016/j.fm.2020.103452.
- 500 Yachison, C.A., Yoshida, C., Robertson, J., Nash, J.H.E., Kruczkiewicz, P., Taboada, E.N.,
501 Walker, M., Reimer, A., Christianson, S., Nichani, A., PulseNet Canada Steering, C.,
502 Nadon, C., 2017. The validation and implications of using whole genome sequencing as
503 a replacement for traditional serotyping for a national *Salmonella* reference laboratory.
504 Front. Microbiol. 8, 1044. <https://doi.org/10.3389/fmicb.2017.01044>.
- 505 Yoshida, C.E., Kruczkiewicz, P., Laing, C.R., Lingohr, E.J., Gannon, V.P., Nash, J.H.,
506 Taboada, E.N., 2016. The *Salmonella* in silico typing Resource (SISTR): an open
507 webaccessible tool for rapidly typing and subtyping draft *Salmonella* genome assemblies.
508 PloS One 11 (1), e0147101. <https://doi.org/10.1371/journal.pone.0147101>.

- 509 Zhang, S., Yin, Y., Jones, M.B., Zhang, Z., Deatherage Kaiser, B.L., Dinsmore, B.A.,
510 Fitzgerald, C., Fields, P.I., Deng, X., 2015. *Salmonella* serotype determination utilizing
511 high-throughput genome sequencing data. *J. Clin. Microbiol.* 2015, 53(5):1685-1692.
- 512 Zhang, S., Den-Bakker, H.C., Li, S., Chen, J., Dinsmore, B.A., Lane, C., Lauer, A.C., Fields,
513 P.I., Deng, X., 2019. SeqSero2: rapid and improved *Salmonella* serotype determination
514 using whole genome sequencing data. *Appl. Environ. Microbiol.* 85(23). Pii, e01746-19.
515 Doi: 10.1128/AEM.01746-19.
- 516 Zhou, Y., Liang, Y., Lynch, K.H., Dennis, J.J., Wishart, D.S., 2011. PHAST: a fast phage
517 search tool. *Nucleic Acids Res.* 39, 347 - 352.

1 **Table 1. Overview of ONT raw reads generated using different sequencing workflows.**

Strain	Sequencing workflow No.	Flowcell	Library preparation kit	Basecalling model	Sequencing time	Mean read length	Mean read quality	Number of reads	Reads length N50	Total bases (Mb)	Coverage	> Q7 (Mb)	> Q12 (Mb)	Percentage of high-quality data
FSL R9-0095	1	R9	Rapid	HAC	2hr	5,799.7	12.1	86,580	9,119	502.1	105	502.1	295.6	58.87%
					24hr	5,816.1	12.1	555,409	8,926	3,230.3	673	3,230.3	1,819.9	56.34%
	2	R9	Rapid	Base modified	2hr	4,450.2	11.7	132,018	8,607	587.5	122	586.3	319.0	54.41%
					24hr	4,441.2	11.7	812,214	8,580	3,607.1	751	3,600.5	1,942.6	53.95%
	3	R9	PCR kit	HAC	2hr	4,431.7	12.5	100,147	5,930	443.8	92	443.8	298.4	67.24%
					72hr	3,945.2	10.7	1,846,360	5,745	7,284.2	1,518	7,283.4	2,213.0	30.38%
	4	R9	1D2	HAC	2hr	12,570.	7.3	84,118	25,490	813.6	170	642.4	38.3	5.96%
					48hr	12,489.	7.2	800,235	25,417	7,649.3	1,594	5,921.2	346.5	5.85%
	5	R10	Rapid	HAC	2hr	9,101.6	9.5	30,805	14,162	280.4	58	280.4	0.0	0%
					24hr	9,237.1	9.3	214,847	14,158	1,983.6	413	1,983.6	0.0	0%
FSL R8-3858	1	R9	Rapid	HAC	2hr	7,026.3	12.0	128,056	12,443	899.7	187	899.6	536.0	59.58%
					48hr	7,081.5	11.5	2,091,542	12,969	14,811.0	3,086	14,809.9	6,799.6	45.91%
	2	R9	Rapid	Base modified	2hr	7,056.8	11.8	140,039	12,840	988.2	206	988.0	527.0	53.34%
					48hr	7,081.5	11.5	2,091,542	12,969	14,811.0	3,086	14,809.9	6,799.6	45.91%
	3	R9	PCR kit	HAC	2hr	4,516.6	12.4	104,149	5,781	470.4	98	470.4	319.1	67.84%
					72hr	4,484.1	11.4	4,424,664	5,746	19,840.7	4,133	19,839.9	8,596.0	43.33%
	4	R9	1D2	HAC	2hr	16,224.	7.4	84,081	30,317	1,098.7	229	884.5	32.6	3.69%
					48hr	15,497.	7.4	643,208	29,007	8,045.2	1,676	6,429.2	298.8	4.65%
	5	R10	Rapid	HAC	2hr	7,894.2	9.2	76,008	13,929	600.0	125	600.0	0.0	0%
					48hr	9,092.8	9.3	471,007	14,312	4,282.8	892	4,282.8	0.0	0%

2 **Table 2. Overview of ONT assembled contigs through data analysis pipeline 1.**

Strain	Sequencing workflow No.	Flowcell	Library Preparation kit	Basecalling model	Sequencing time	Total length (bp)	Contigs Numbers	N50 Length (bp)	SNPs	InDels	Matching identity (%)	High quality SNPs
FSL R9-0095	1	R9	Rapid	HAC	2hr	4,792,171	3	4,736,332	4,778	3,116	99.83	4
					24hr	4,789,608	2	4,741,208	4,492	2,534	99.85	2
	2	R9	Rapid	Base modified	2hr	4,825,067	4	4,740,589	31	901	99.98	0
					24hr	4,741,424	1	4,741,424	29	652	99.98	0
	3	R9	PCR kit	HAC	2hr	4,667,484	49	183,328	42	493	99.99	0
					72hr	4,561,879	48	174,019	245	1,001	99.97	0
	4	R9	1D2	HAC	2hr	4,742,094	1	4,742,094	5,043	5,840	99.77	1
					48hr	4,790,286	2	4,741,889	5,200	4,911	99.78	4
	5	R10	Rapid	HAC	2hr	4,790,793	2	4,742,687	619	1,672	99.95	1
					24hr	4,783,849	2	4,735,698	811	2,179	99.93	1
FSL R8-3858	1	R9	Rapid	HAC	2hr	4,665,249	1	4,665,249	4,831	2,617	99.84	100
					48hr	4,665,134	1	4,665,134	4,854	2,705	99.83	73
	2	R9	Rapid	Base modified	2hr	4,666,028	1	4,666,028	33	677	99.98	1
					48hr	4,666,136	1	4,666,136	38	666	99.98	1
	3	R9	PCR kit	HAC	2hr	4,583,190	31	273,714	56	483	99.98	1
					72hr	4,500,549	33	180,233	59	407	99.99	1
	4	R9	1D2	HAC	2hr	4,666,639	1	4,666,639	5,605	6,524	99.74	69
					48hr	4,698,102	2	4,665,652	5,964	6,724	99.72	61
	5	R10	Rapid	HAC	2hr	4,666,923	1	4,666,923	691	1,919	99.94	4
					48hr	4,651,530	2	3,508,306	608	1,873	99.95	3

3

4 **Table 3. Result of predicted prophage for strain FSL R8-3858.**

Sequencing technology	Sequencing workflow No.	Flowcell	Library preparation kit	Basecalling model	Sequencing time	Prophage □ ₁₅			Prophage □ ₃₄		
						Length (kb)	Completeness	Score	Length (bp)	Completeness	Score
Illumina	1	R9	Rapid	HAC	2hr	25.6	Incomplete	20	36.8	Questionable	80
					48hr	45.9	Intact	130	63.7	Questionable	90
ONT	2	R9	Rapid	Base modified	2hr	44.3	Intact	100	43.1	Questionable	90
					48hr	46	Intact	100	43.7	Questionable	80
	3	R9	PCR kit	HAC	2hr	44.3	Intact	100	44.1	Questionable	90
					72hr	44.3	Intact	100	45.1	Questionable	90
	4	R9	1D2	HAC	2hr	45.7	Intact	150	43.7	Questionable	90
					48hr	45.6	Intact	150	33.2	Incomplete	30
5	R10	Rapid	HAC	2hr	46	Intact	110	60.1	Questionable	90	
				48hr	44.3	Intact	130	47.1	Questionable	90	

5 **Note:** When the score was above 90, the completeness was recorded as intact; when the score was
6 between 70 and 90, the completeness was recorded as questionable; when the score was below 70, the
7 completeness was recorded as incomplete. More information about the result interpretation for scoring
8 the prophage region can be found on PASTER API (<http://phaster.ca>).

9 **Figure legends**

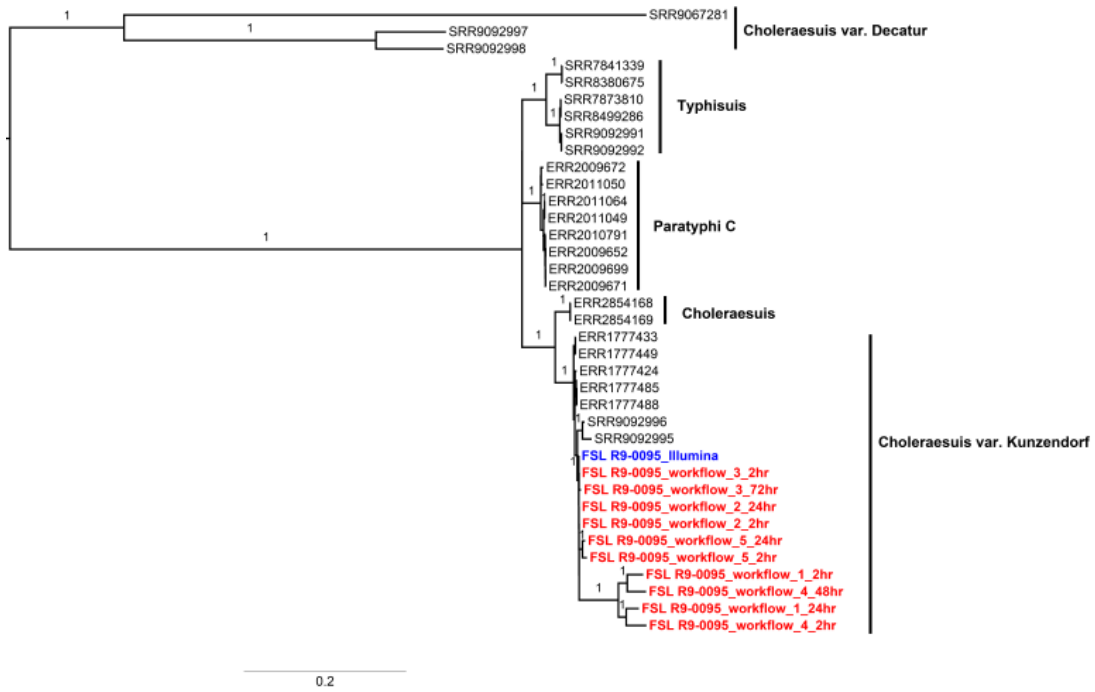
10 **Figure 1. Maximum parsimony tree of the serotype formula 6,7:c:1,5 based on k-mer-based SNP**

11 **analysis.** The tree was built using kSNP3 with the core SNPs identified among both Illumina and
12 ONT sequences of the strain FSL R9-0095 and 26 sequences downloaded from the NCBI. Clades
13 of Typhisuis, Paratyphi C, Choleraesuis, Choleraesuis var. Kunzendorf, Choleraesuis var. Decatur
14 were annotated separately. Only high support in the analysis (bootstrap \geq 90%) are labeled in
15 this tree. The tree is midpoint rooted and the scale axis is provided below the tree.

16 **Figure 2. Maximum likelihood phylogenetic tree of the serotype formula 6,7:c:1,5 using the**

17 **CFSAN hqSNP pipeline.** The tree was constructed with PhyML using hqSNPs identified among
18 Illumina and ONT sequences of the strain FSL R9-0095 and 26 sequences downloaded from the
19 NCBI. Illumina sequences of FSL R9-0095 was used as the reference genome. (a) The tree
20 includes both the Illumina and the ONT sequencing data of strain FSL R9-0095. Clades of
21 Typhisuis, Paratyphi C, Choleraesuis, Choleraesuis var. Kunzendorf, Choleraesuis var. Decatur
22 are annotated separatel. (b) The data of strain FSL R9-0095 replaced with SNPs from two hrs
23 ONT sequencing data using sequencing workflow 2. (c) The tree only includes Illumina
24 sequencing data of strain FSL R9-0095, which were used as a benchmark phylogenetic tree. Only
25 high support in the analysis (bootstrap \geq 90%) are labeled in this tree. The tree is midpoint
26 rooted and the scale axis is provided below the tree.

27
28

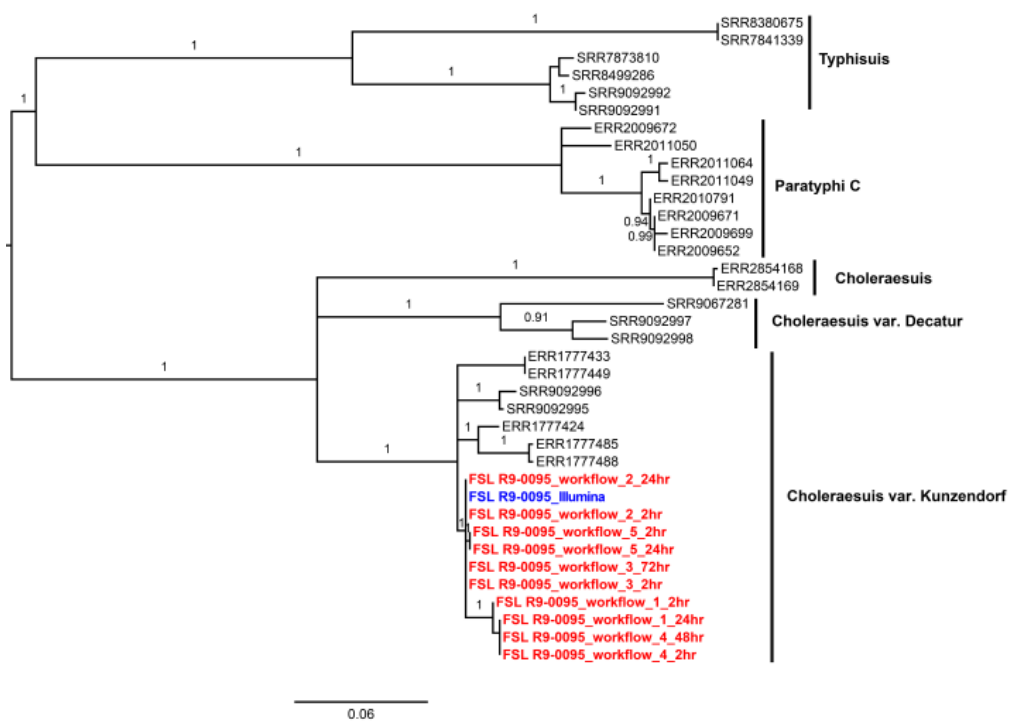


29

30

Figure 1

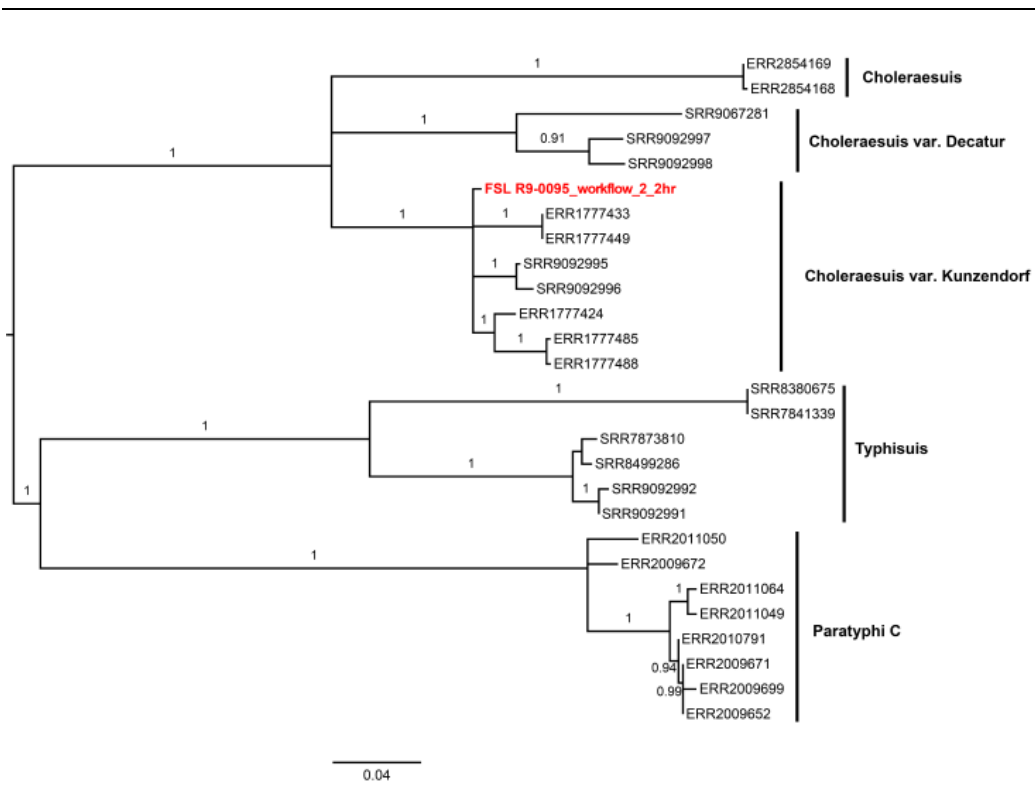
31



32

33

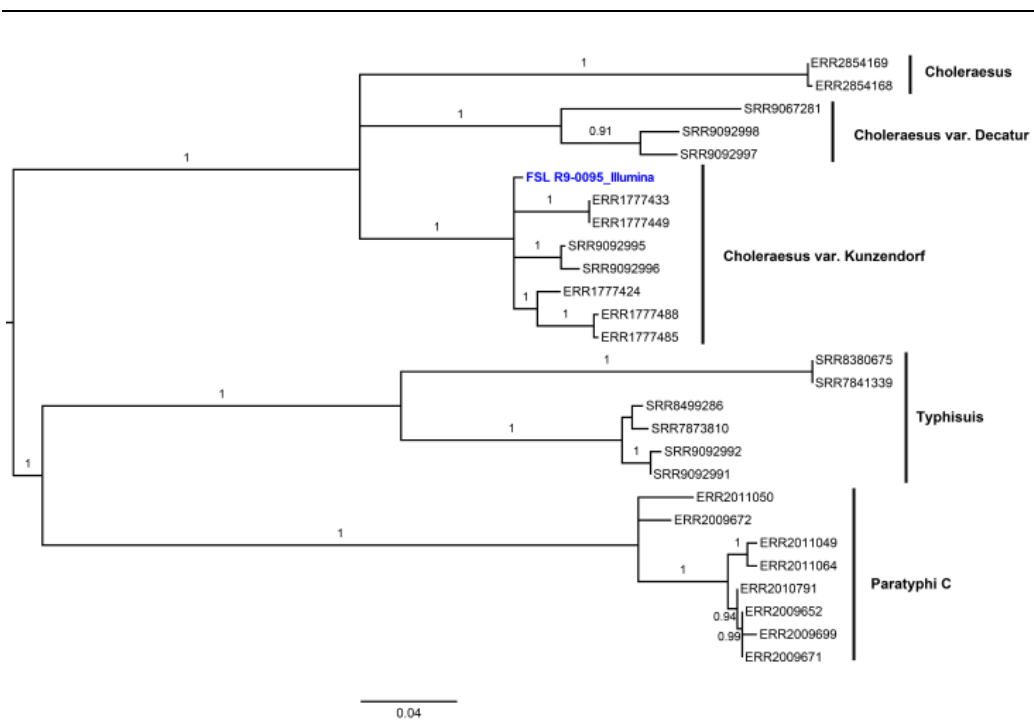
Figure 2a



34

35

Figure 2b



36

37

Figure 3c