# gapFinisher: a reliable gap filling pipeline for SSPACE-LongRead scaffolder output

**Juhana I Kammonen** [Corresp., 1] , **Olli-Pekka Smolander** [1] , **Lars Paulin** [1] , **Pedro AB Pereira** [1] , **Pia Laine** [1] , **Patrik Koskinen** [2] , **Jukka Jernvall** [1] , **Petri Auvinen** [1]

[1] Institute of Biotechnology, University of Helsinki, Helsinki, Finland

[2] Blueprint Genetics Ltd., Helsinki, Finland

Corresponding Author: Juhana I Kammonen
Email address: juhana.kammonen@helsinki.fi

Unknown sequences, or gaps, are largely present in most published genomes across public databases. Gap filling is an important finishing step in *de novo* genome assembly, especially in large genomes. The gap filling problem is nontrivial and while many computational tools exist partially solving the problem, several have shortcomings as to the reliability and correctness of the output, i.e. the gap filled draft genome. SSPACE-LongRead is a scaffolding software that utilizes long reads from multiple third-generation sequencing platforms in finding links between contigs and combining them. The long reads potentially contain sequence information to fill the gaps, but SSPACE-LongRead currently lacks this functionality. We present an automated pipeline called gapFinisher to process SSPACE-LongRead output to fill gaps after the actual scaffolding. gapFinisher is based on controlled use of a gap filling tool called FGAP and works on all standard Linux/UNIX command lines. We conclude that performing the workflows of SSPACE-LongRead and gapFinisher enables users to fill gaps reliably. There is no need for further scrutiny of the existing sequencing data after performing the analysis.

1 Helsinki, Finland 07 Dec 2017       MANUSCRIPT TEXT

2 **gapFinisher: a reliable gap filling pipeline for SSPACE-LongRead scaffolder output**

3 Juhana I Kammonen (1)*, Olli-Pekka Smolander (1), Lars Paulin (1), Pedro AB Pereira (1), Pia Laine (1),

4 Patrik Koskinen (1,2), Jukka Jernvall (3) & Petri Auvinen (1)

5 (1) DNA Sequencing and Genomics Laboratory, Institute of Biotechnology, University of Helsinki,

6 Finland

7 (2) Blueprint Genetics Ltd. Finland

8 (3) Evolutionary Phenomics Group, Institute of Biotechnology, University of Helsinki, Finland

9 *corresponding author ( juhana.kammonen@helsinki.fi )

10 BIOGRAPHICAL NOTE

11 The corresponding author is a PhD student representing the Saimaa ringed seal genome project

12 (http://www.saimaaringedseal.org). The PhD thesis topic covers the full genome sequencing and genome

13 annotation of the Saimaa ringed seal (*Pusa hispida saimensis*), considered one of the national animals of

14 Finland and one of the most endangered animals in the world.

15 KEYWORDS

16 genome assembly; draft genomes; scaffolding; next-generation sequencing; long read technologies

17 ABSTRACT

18   Unknown sequences, or gaps, are largely present in most published genomes across public databases. Gap

19   filling is an important finishing step in *de novo* genome assembly, especially in large genomes. The gap

20   filling problem is nontrivial and while many computational tools exist partially solving the problem,

21   several have shortcomings as to the reliability and correctness of the output, i.e. the gap filled draft

22   genome. SSPACE-LongRead is a scaffolding software that utilizes long reads from multiple third-

23   generation sequencing platforms in finding links between contigs and combining them. The long reads

24   potentially contain sequence information to fill the gaps, but SSPACE-LongRead currently lacks this

25   functionality. We present an automated pipeline called gapFinisher to process SSPACE-LongRead output

26   to fill gaps after the actual scaffolding. gapFinisher is based on controlled use of a gap filling tool called

27   FGAP and works on all standard Linux/UNIX command lines. We conclude that performing the

28   workflows of SSPACE-LongRead and gapFinisher enables users to fill gaps reliably. There is no need for

29   further scrutiny of the existing sequencing data after performing the analysis.

30   INTRODUCTION

31   Gap filling is one of the final phases of *de novo* genome assembly. First, assembly algorithms produce

32   contiguous sequences of overlapping sequencing reads, commonly known as contigs. A contig is a

33   continuous DNA sequence entity without any ambiguities, unknown bases marked as N. Second, the

34   contigs are connected into longer fragments, scaffolds, using specialized sequencing read data. Until the

35   development of long read technologies, the data used to be primarily mate-pair reads, known also as

36   jumping reads. The mate-pair libraries are usually made of size selected DNA fragments, where fragment

37   size is usually in the order of thousands of base pairs. The ends of these fragments are sequenced and

38   resulting reads are used for scaffolding. Currently long continuous reads e.g. from Pacific Biosciences

39   (Menlo Park, California, United States) RS II or Sequel third-generation sequencing platforms are

40   commonly used. While the scaffolding step links and orders the contigs, it usually leaves variable amounts

41   of unknown sequences, strings of N-characters, in between them. These unknown sequences are called

42  gaps. Finally, the gap filling stage aims to resolve these unknown sequences with additional sequencing

43  data, (e.g. Boetzer *et al.*, 2011; Boetzer & Pirovano, 2012) or without additional data (e.g. Li & Copley,

44  2013).

45  In this paper, we present an automated gap filling pipeline called gapFinisher. We pursue a solution to the

46  gap filling problem that utilizes long reads and only unaltered draft genomes to prevent any loss of data.

47  We set strict alignment parameters for the gap filling stage to ensure correctness and uniqueness of the

48  filled gaps. We conclude that applying gapFinisher enables efficient and reliable gap filling by controlling

49  the use of the FGAP algorithm (Piro *et al*., 2014). Furthermore, gapFinisher time complexity proves linear

50  with respect to the size of the input. The system requirements are MATLAB Compilation Runtime (MCR)

51  for FGAP and a Perl interpreter for SSPACE-LR. Besides these, the gapFinisher pipeline does not require

52  any additional software to be installed.

53  **Gap filling**

54  SSPACE-Standard (Boetzer *et al*., 2011) and SSPACE-LongRead (SSPACE-LR) (Boetzer & Pirovano,

55  2014) are scaffolding tools for paired-end (also mate-pair) reads and long continuous reads, respectively.

56  While this software is available free for academic users, both are commercial products and upgrades and

57  most of the support require a proprietary license. SSPACE-Standard is commonly applied in the first

58  scaffolding step where contigs are oriented and ordered into longer connected sequences and it accepts

59  paired-end data from any next-generation sequencing technology if read-orientation information and mean

60  values and standard deviations of the insert sizes for each read library are provided. SSPACE-LR utilizes

61  Pacific Biosciences filtered subreads (CLR = *Continuous Long Reads*) in finding long links between

62  contigs or existing scaffolds and combining them into "superscaffolds" with new gaps introduced between

63  the sequences. SSPACE-LR first maps the long reads into the contig assembly using the BLASR aligner

64  specialized for long read alignment (Chaisson & Tessler, 2012). Based on these alignments, contigs are

65  then linked into scaffolds and N-characters (gaps) are placed between the connected contigs. While the

66  CLR reads contain information of the actual nucleotide sequence in the gaps, this feature is not exploited

67    in the current version of SSPACE-LR (version 1.1). However, the software can report the exact

68    information about which reads were associated when creating the new scaffold and the new gap(s). In the

69    gapFinisher pipeline, we actually utilize this information to fill the gaps in the newly created scaffolds on

70    the go.

71    A central part of gap filling is the alignment of long sequences against the contigs. This is challenging due

72    to the relatively high error-rates of contemporary long read data (Laver *et al*., 2015) and the sequencing

73    errors, e.g. (Nakamura *et al*., 2011; Schirmer *et al*., 2015) and local misassemblies at the contig level

74    (Boetzer & Pirovano, 2014). The BLAST local alignment tool (Altschul *et al*., 1990) is the most

75    commonly used approach for the identification of areas of high similarity between multiple sequences.

76    Different scaffolding and gap filling tools apply BLAST either directly (Piro *et al*., 2014), or the method is

77    refined (Chaisson & Tessler, 2012) and applied (English *et al*., 2012; Boetzer & Pirovano, 2014). All tools

78    based on BLAST contain multiple parameters, e.g. for mismatches and gaps, affecting their ability to

79    detect non-perfect matches and it is not always clear how these should be defined.

80    Several gap filling software exist. GapFiller (Boetzer & Pirovano, 2012) is a commercial program by the

81    authors of SSPACE-tools and often used together with them. GapFiller uses paired-end read information

82    to fill in sequences at contig ends where overlapping reads reach into the gap created on the SSPACE-

83    Standard step by mate-pair reads. Where reads are unable to span the whole length of the unknown

84    sequence, the gap is not completely filled and unknown bases (N-characters) will remain in the output

85    version of the draft genome (Boetzer *et al*., 2011).

86    PBJelly (English *et al*., 2012) is a scaffolding and gap filling tool integrated into the Pacific Biosciences

87    (PacBio) SMRT Analysis software suite, the main user interface for data analysis using PacBio long reads.

88    In comparison to other gap filling tools, the PBJelly is run in six separate stages (setup, mapping, support,

89    extraction, assembly and output) and requires additional software libraries, most notably the SMRT Portal

90    software suite and BLASR (Chaisson & Tessler, 2012). Despite it is possible to run PBJelly in a single-

91    core computer, the workflow is clearly designed for high-throughput computing in a grid, e.g. the Sun

92    Grid Engine (Gentsch, 2001). A peculiar default feature of PBJelly is that it by default inflates short gaps

93    (< 25 bp) to a length of exactly 25 bp with the apparent purpose of emphasizing the location of the gaps

94    (English *et al*., 2012).

95    Gap2Seq (Salmela *et al*., 2016) provides a purely computational solution to the gap filling problem. It

96    works well on prokaryote genomes but does not scale well to larger genomes, where repetitive sequences

97    confuse the algorithm and the sheer size of the genome makes running times infeasibly long.

98    FGAP (Piro *et al*., 2014) is a gap filling tool that utilizes various types of read data and BLAST

99    alignments to find and fill gaps in draft genomes. The BLAST utility is bundled with the release version of

100   FGAP, but a MATLAB Compilation Runtime is required. The gapFinisher pipeline presented in this paper

101   is based on FGAP and enables more reliable and controlled gap filling.

102   **State of the art**

103   Although FGAP efficiently reduces the number of gaps in various draft genomes (Piro *et al*., 2014), the

104   tool has the rather troubling feature of setting no limit to the number of times an input read is used in gap

105   filling should the BLAST alignment return multiple good hits (Fig. 1). With the default setting of FGAP,

106   undesired multiple alignments of query sequences may occur due to repetitive regions in the draft genome,

107   or overly lenient alignment parameters for the ends of the query sequences (Fig. 1). We could verify this

108   behaviour on an FGAP test run with a preliminary draft genome of an unpublished marine mammal from

109   the *Phocidae* family (Table S4). Ideally, gap filling should be a unique process in the sense that a single

110   input long read would find a single good alignment in the draft genome and fill any gaps in that single

111   location.

112   (Figure 1 here)

113   Repeat masking may improve the scaffolding and gap filling of highly repetitive draft genomes. For

114   example, it has been estimated that more than 60% of the 3,3 Gb modern human (*H. sapiens*) genome

115   consists of repetitive sequences, (e.g. de Konig *et al*., 2011). With the repetitive sequences at contig ends

116   eliminated, the scaffolding / gap filling algorithms are less likely to make incorrect alignments. One

117  example of repeat masking software tools is RepeatMasker (Smit *et al.*, 2013) which finds short and long

118  interspersed elements as well as simple repeats in the input genomic sequence. RepeatMasker may mask

119  coding regions of the input genome, especially those located at the terminal regions of open reading

120  frames (ORFs). Furthermore, RepeatMasker may mask some shorter potential element-coding sequences

121  such as ribosomal RNAs (Smit *et al.*, 2013). Repeat masking may lower the inherent risk of incorrect

122  alignments or multiple alignments in the contig ends. In this paper, however, we pursue a solution that

123  utilizes only unaltered (unmasked) draft genomes to prevent any loss of data.

124  Solving short gaps of e.g. 1-20 base pairs in length by simple read alignment maps produced by e.g. the

125  Burrows-Wheeler Aligner (Li & Durbin, 2009) or the Bowtie 2 aligner (Langmead & Salzberg, 2012) is

126  not investigated in detail in this study, but may be one of the prospects of solving the gap filling problem

127  for short gaps. For instance, some singular unknown bases and short N-sequences at gap edges are solved

128  by the re-assembly stage of the Pilon assembly polishing tool, where an alignment map file can be

129  supplied as input and a specific option set for gap filling (Walker *et al.*, 2014).

130  MATERIALS & METHODS

131  The current release of gapFinisher works only on the output of SSPACE-LongRead (Boetzer & Pirovano,

132  2014). The basic workflow of gapFinisher is illustrated in Figure 1c and in further detail in Figure 2. Thus,

133  before running gapFinisher, the user must successfully run SSPACE-LongRead for a dataset at least once.

134  It is imperative to have the "-k" option enabled when running SSPACE-LongRead. This setting will create

135  the critical "inner-scaffold-sequences" subdirectory that contains for each superscaffold the references to

136  the actual long read sequences (one or more) that created the scaffold. The pipeline will not run if this

137  directory does not exist. When successful, gapFinisher then works as follows (*cf.* Figure 2):

138      1.   Index the draft genome FASTA file and the long read FASTA file

139      2.   Generate a list of names of all superscaffolds SSPACE-LongRead (-k 1 option enabled) has

140           created
141      3.   For each superscaffold in the list:
142               a.   Create a new FGAP working directory for the current superscaffold
143               b.   Fetch all full CLR reads associated with the current superscaffold
144               c.   For each of the CLR reads associated with the current superscaffold:
145                        i.   Execute FGAP using the current superscaffold as draft and the CLR read as input
146                        ii.  If FGAP could fill (one or more) gaps in the current superscaffold, save FGAP

147                             output as the new draft for the current superscaffold

148      4.   Compile results from each working directory as filled_scaffolds.fasta
149      5.   Compile filled_scaffolds.fasta and the unfilled/untouched scaffolds from the original draft genome

150           as scaffolds_gapfilled_FINAL.fasta
151      6.   [optional] Clean the working directories (to save disk space).

152      The rapid fetching of reads is based on the operation of the fastaindex (step 1 above) and fastafetch (step

153      2b above) utilities of the exonerate toolkit (Slater & Birney, 2004) v. 2.4.0. Precompiled executables of

154      these utilities are bundled with the gapFinisher release and fully integrated into the workflow of the

155      gapFinisher pipeline.

156      (Figure 2 here)

157      When using Pacific Biosciences filtered subreads with SSPACE-LongRead, it is in theory possible, that

158      separate reads originating from the same well of the PacBio SMRT cell are aligned into separate places by

159      the BLASR aligner (cf. Fig. 1a and Fig. 1b). Filtered subreads from the same well of the SMRT cell

160      always originate from the same molecule and thus should align to locations close to one another. The

161      legacy BLASR (Chaisson & Tessler, 2012) version that SSPACE-LongRead is using has no formal

162      assertion for this. Hence, we set gapFinisher to keep track of the origins of the filtered subreads. The

163      pipeline issues an appropriate warning when gap filling under conflicting read origin is about to happen

164      and aborts the filling of the gap in question. Conflicting read origins further indicate potential errors in the

165     scaffolding step. Consequently, the location and read information of the conflict are included in the

166     warning message and logged.

167     (Table 1 here)

168     Here, we predisposed seven separate sequencing read datasets from both bacterial and eukaryote

169     organisms (Table 1) to *de novo* assembly and scaffolding. Finally, we performed gap filling on the created

170     scaffolds with gapFinisher (Table 2). First, we had two *Escherichia coli* (*E. coli*) bacterial genome drafts.

171     Second, we extended the analysis to a set of further four bacterial genomes: *Bibersteinia trehalosi,*

172     *Mannheimia haemolytica, Francisella tularensis* and *Salmonella enterica*. The bacterial read data are the

173     same that were used as test data for the SSPACE-LongRead scaffolder (Boetzer & Pirovano, 2014) and are

174     available at: http://www.cbcb.umd.edu/software/PBcR/closure/index.html and the Sequencing Read

175     Archive (SRA) links therein. For *B. trehalosi*, we used the reference sequence *Bibersteinia trehalosi*

176     USDA-ARS-USMARC-188 (Harhay *et al*., 2014) Since the publication of SSPACE-LongRead, a

177     reference genome has become available to *M. haemolytica* as well (Eidam *et al*., 2013). Finally, we

178     included an unpublished marine mammal (*Phocidae* family) draft genome in final stage with 236,592

179     contigs scaffolded into 10,143 superscaffolds with gaps to get a reference on how gapFinisher performs on

180     a much larger genome. The raw sequencing coverage of the draft genome was on average 25X for the

181     Illumina (San Diego, California, United States) reads and 50X for the PacBio CLR reads (Table 1). When

182     assembled with the miniasm (Li, 2016) using all the PacBio reads, we got an additional "PacBio-only"

183     assembled version of the draft genome with 1,314 contigs which we then scaffolded and gap filled (Table

184     2).

185     For the Illumina MiSeq reads, we further applied the Fast Length Adjustment of SHort Reads (FLASH)

186     protocol that finds overlaps at the ends of the MiSeq paired-end reads and joins the reads if found (Magoč

187     & Salzberg, 2011). Thus, about half of the reads in each MiSeq dataset could be combined to longer initial

188     fragments. This feature is likely to improve the *de novo* genome assemblies while longer initial read

189     length may be enough to span short repeats and indels. The uncombined reads were supplied as additional

190  paired-end libraries in all assemblies. The Roche 454 Genome Sequencer data available for the draft

191  genomes was not utilized here, as our benchmark did not include a suitable assembler, e.g. Newbler

192  (Margulies *et al*., 2005) for these data. Furthermore, the performance of Newbler was evaluated rather

193  extensively in the SSPACE-LongRead original publication (Boetzer & Pirovano, 2014) and in most of

194  cases Newbler could not perform as well as the other short read assemblers.

195  We assembled the draft genomes with the SPAdes (Bankevich *et al*., 2012) and miniasm (Li, 2016)

196  assemblers. SPAdes can employ both Illumina short reads and PacBio CLR reads. In contrast, miniasm

197  only works properly with PacBio CLR reads or other long reads with a sufficient sequencing coverage.

198  This is because the read trimming phase of miniasm requires a read-to-read mapping length of at least

199  2,000 bp with a minimum of 100 bp non-redundant bases (Li, 2016). This condition is not met by the short

200  read datasets used in this study. An additional and a highly useful feature of miniasm is the minidot plot

201  drawing utility and it was used to create the dotplots for comparisons to the reference genomes (Figs. 3

202  and S1).

203  The scaffolding step included the combined use of SSPACE-LongRead (academic license, software

204  version 1.1) (Boetzer & Pirovano, 2014) and the gapFinisher pipeline. We first executed SSPACE-

205  LongRead for all samples to create the superscaffold assemblies for the six bacterial genomes and the

206  unpublished mammal draft genome (*Phocidae* family). The same long read data was applied for the

207  scaffolding of both SPAdes and miniasm contig assemblies. For each scaffold assembly, we then executed

208  gapFinisher to fill the gaps introduced by the scaffolding step. Due to the large size (~2,5 gigabases) of the

209  unpublished mammal genome, the SSPACE-LongRead and gap filling stage for the miniasm assembly had

210  to be executed in two consecutive runs with 25X (50 % of the total coverage) PacBio reads applied to each

211  part. On the other hand, the scaffolding of the mammal SPAdes assembly was executed in five separate

212  stages as part of the actual genome project of the mammal. About 10X coverage of PacBio reads of insert

213  were applied at each stage and gapFinisher executed between the stages. Thus, the results for this

214    assembly (Table 2) show statistics for the final stage and average CLR reads per scaffold is the average of

215    all five stages.

216    We visualized the different stages of the draft assemblies for all genomes using the subplot utility of the

217    MATLAB toolkit (Figs. 3 and S1). Furthermore, we visualized the final stages of the assembly and

218    scaffolding by aligning the reference genomes and the two drafts from the SPAdes and miniasm assembly

219    pipelines with the progressiveMauve algorithm of the Mauve (Darling *et al*., 2004) alignment and

220    visualization tool (Figs. 4 and S2). Mauve reveals the number and similarity of Locally Collinear Blocks

221    (LCBs) between the input sequences.

222    In order to assess the performance of the software, all of the SPAdes, miniasm, SSPACE-LongRead and

223    gapFinisher runs were executed in two separate 64-bit Linux computer environments. First, the bacterial

224    genomes were assembled, scaffolded and gap filled in a single-processor computer running Ubuntu Linux

225    14.04 with 20 GB of RAM, the equivalent to a modern office workstation with a small RAM extension.

226    The processor was Intel (Santa Clara, California, United States) Core (TM) with a frequency of 3,2 GHz.

227    Second, we built the mammal genome in a multi-core supercomputer running Ubuntu Linux 14.04 with 1

228    TB of RAM and using 16 Advanced Micro Devices (Sunnyvale, California, United States) Opteron(TM)

229    processors with a frequency of 2,5 GHz each. The latter setup is equivalent to a small-scale computer

230    cluster.

231    RESULTS

232    The results are presented both from the viewpoint of how finished the draft genomes are before and after

233    gapFinisher and how gapFinisher performs in general and with respect to FGAP. Key statistics of the

234  assembly benchmark results were compiled (Fig. 5) and the genome alignments to the bacterial reference

235  genomes were visualized (Figure .3 and Figure. 4 and supplementary figures S1 and S2).

## Genomes

237  Regarding *de novo* assembly of the genomes, we noticed similar behaviour of the SPAdes assembler as

238  reported by the authors of the SSPACE-LongRead (Boetzer & Pirovano, 2014). Namely, that the SPAdes

239  assembly pipeline introduced repeats at the ends of the contigs that evidently prohibit many CLR reads

240  from aligning into the contig ends and thus the scaffold assembly is left with a higher number of

241  uncombined sequences than expected (Table 2 and Fig. 4a). Nevertheless, scaffolding with SSPACE-

242  LongRead reduced the number of total sequences in all the assemblies. This was especially evident in the

243  *Mannheimia haemolytica* draft genome, where SSPACE-LongRead reduced the number of sequences in

244  the draft assembly from 112 to 17 (84,8 % reduction). A notable increase in basic assembly statistics, such

245  as the N50 contig length and number of sequences, was observed throughout (Table 2). The miniasm

246  assembler (Li, 2016) outperformed the assemblers used in the SSPACE-LongRead test assemblies

247  (Boetzer & Pirovano, 2014) and the SPAdes assembler (Bankevich *et al.*, 2012) in our benchmark in terms

248  of number of output contigs, N50 and gap length (Table 2). On the other hand, the median similarity of the

249  alignments to the bacterial reference genomes is lower across all bacterial draft genomes from the

250  miniasm pipeline (Figs. 4 and S2).

251  It is evident from the assembly results that both the SPAdes and miniasm assemblers are optimized for the

252  *E. coli* K12 genome: The number of SPAdes assembly contigs was the lowest of the bacterial assemblies

253  in this study, namely 35 (Table 2). The miniasm assembly of the *E. coli* K12 genome was a single

254  sequence (Fig. S1 and Table 2) and thus was the only draft genome not to require scaffolding or gap

255  filling. Furthermore, miniasm was able to construct the full *E. coli* K12 genome from PacBio reads in 3

256   minutes (Fig. 5 and Table S3). The final assembly consists of a single long bacterial genome (Table 2) in 4

257   Locally Collinear Blocks (LCB's) according to progressiveMauve (Darling *et al*., 2004) alignment (Table

258   2 and Fig. S2a). The contig assembly results for the other bacterial genomes were more variable with both

259   SPAdes and miniasm (Table 2 and Figs. 1 and 2).

260   Regarding the overall output of the assemblers, miniasm consistently reports zero N's at the contig

261   assembly stage (Table 2). Furthermore, the miniasm contig assemblies are more contiguous in the sense

262   that they consist of less sequences when compared to the SPAdes assemblies in all cases (Table 2).

263   However, the SPAdes assemblies report some gapped sequences already at the contig assembly of *E. coli*

264   *O157* (3 bp), *B. trehalosi* (2 bp), *M. haemolytica* (35 bp) and *S. enterica* (655 bp) (Table 2).

265   Evidently, gapFinisher can fill about 50 % of the gapped sequence (Table 2) in the scaffolded draft

266   genomes and retains the structure of the genomes in all cases (Figs. 3 and 4 and Supplementary figures S1

267   and S2). The lowest percentage of gaps filled was with the second stage of the mammal genome miniasm

268   scaffolding (4,1 %) and the highest percentage of gaps filled was with the scaffolding of the B. trehalosi

269   SPAdes assembly (85,7 %). At the nucleotide level, several kilobases of gapped sequence is being filled in

270   all cases (Table 2). No large insertions, deletions or inversions are introduced by the gap filling stage with

271   gapFinisher (Table 2 and Figs. 3 and S1). The bacterial initial assemblies were refined to scaffolds using

272   PacBio filtered subreads. There were no cases of gapFinisher warning about separate reads from same

273   SMRT cell well attempting to fill disparate gaps in any of the bacterial genomes.

274   (Figures 3 and 4 here) (Table 2 here)

275   **Performance**

276   (Figure 5 here)

277    Besides MATLAB Compilation Runtime and a Perl (Christiansen *et al.*, 2012) interpreter, gapFinisher

278    does not require any other software to be installed. Furthermore, the gapFinisher pipeline is contained in a

279    single phase, namely the actual execution of the gap filling, where e.g. the PBJelly (English *et al.*, 2012)

280    pipeline has six separate phases.

281    Due to the serial design of the pipeline, gapFinisher running time holds quite neatly at about 3-5 wall-

282    clock seconds per CLR read per scaffold (Table S3). Thus, gapFinisher time complexity can be regarded

283    as linear with relation to the number of input scaffolds and the total coverage of the scaffolding reads.

284    Where average number of CLR reads per created scaffold was high, as was the case with the bacterial

285    genomes of *F. tularensis, M. haemolytica* and *S. enterica*, gapFinisher running time in single-core mode

286    was notably higher (Fig. 5 and Table S3).

287    We studied random access memory (RAM) use of gapFinisher (Fig. 5 and Table S3). We used a built-in

288    Linux utility (/usr/bin/time) to measure the peak RAM use during each of the assembly stages. Again, the

289    serial design of gapFinisher keeps the RAM use of the gap filling stage at all but nominal level (Fig. 5a),

290    also in the case of a much larger genome (Fig. 5c). In general, the peak RAM use of less than 1 GB we

291    detected in all cases means that gapFinisher could be executed in almost any Linux computer, even most

292    tablets. Nevertheless, the preceding assembly steps tend to use significantly more RAM (Fig. 3a and Fig.

293    3c). In particular, the larger mammal genome used more than 500 GB of RAM in the contig assembly

294    stage and more than 80 GB of RAM in the SSPACE-LongRead stage (Table 2).

295    DISCUSSION

296    Gap filling is a non-trivial problem with many existing solutions today in the form of software tools. The

297    correctness of the outputs of different tools is variable. For a large genome under assembly, the default

298    parameter settings of FGAP clearly are too lenient and may lead to incorrect gap filling in large draft

299 genomes (Table S4). Repeat masking before gap filling with FGAP alone may be recommended,

300 especially because FGAP utilizes BLAST (Altschul *et al*., 1990) directly for the long read alignment.

301 Typically, contig assemblies do not contain any unknown sequence (N-characters) and the output of

302 miniasm correctly follows this principle (Table 2). However, it is evident from the SPAdes assembler

303 results that a small number of N's may be introduced already at the contig assembly stage (Table 2). This

304 may be due to the N's present in the actual read data that is not uncommon for Illumina sequencing reads

305 but is more unusual for PacBio reads.

306 gapFinisher is not able to fill all gapped sequences in the draft assembly (Table 2). This is because the

307 CLR reads of the Pacific Biosciences platform do contain base-call errors, (e.g. Laver *et al*., 2015) and

308 gapFinisher employs a strict alignment scheme of the long reads and only fills a gap when a reasonably

309 correct alignment of known sequences at the gap edges is found (Figs. 1c and 2). Consequently, it is

310 possible that some gaps are prevented from filling despite the evidence being there. A solution is to run

311 gapFinisher on less strict parameters and then confirm the correctness of the result using other alignment

312 tools. Nevertheless, gapFinisher with the default settings can reduce the amount of gapped sequence in the

313 example draft genomes by about 50 % in general (Table 2).

314 Regarding the use of filtered subreads in the bacterial genome assemblies of this study, gapFinisher did

315 not detect any cases where separate reads from the same SMRT cell well would have filled disparate gaps

316 in the genomes. In applications where conflicting read origins could be a problem, it can be circumvented

317 by producing reads of insert from the filtered subreads, albeit with the expense of genome level coverage.

318 On the other hand, the reads of insert pipeline improves the overall quality of the reads which leads to

319 more reliable alignments. Checking the read origin of the filtered subreads is a valuable additional

320 correctness feature of the gapFinisher pipeline not present in the other gap filling tools presented in this

321 study.

322 We found that the time complexity of gapFinisher is approximately linear with respect to the number of

323 input scaffolds and long read evidence related to each of the scaffolds (Fig. 5 and Supplementary table

324 S3). While the peak RAM use of gapFinisher stays at a nominal level in all the cases of small and large

325 genomes (Fig 5a and Fig. 5c), the runtime varies significantly, even in small genome assemblies (Fig 5b).

326 This feature may be optimized in future development versions. If the user can run gapFinisher in a

327 supercomputer cluster, it is possible to specify the number of threads (option -t) and the utility will divide

328 the input scaffolds into even parts, splitting the total running time by the number of processors assigned.

329 In the case of our datasets, the parallelization would have significantly reduced the runtime of gapFinisher

330 in the gap filling of bacterial genomes *M. haemolytica* and *S. enterica* (Fig. 5b and Supplementary table

331 S3) and the effects could be clearly seen in the case of the mammal genome gap filling with 16 processors

332 in use and parallelizing the workflow (Fig. 5c). gapFinisher is designed to work on all standard

333 Linux/UNIX distributions on command line with as little dependencies as possible. Aside for having to

334 first perform the actual scaffolding using SSPACE-LongRead, all the user needs to do is download

335 gapFinisher and run it.

336 No matter which next-generation sequencing platform is in use, there exists a distinct base-call error

337 profile affecting the output and the quality of the sequenced reads. Previously, sequence-specific

338 systematic miscalls have been reported in the output of Illumina Genome Analyzer II platform (Dohm *et*

339 *al*., 2008; Nakamura *et al*., 2011). Evidently, the more recent Illumina MiSeq platform is affected by the

340 same miscall profile to some extent (Kammonen *et al*., 2015; Schirmer *et al*., 2015). The presence of a

341 relatively high error-rate can also not be disputed in current high-throughput sequencing of long reads

342 (Laver *et al*., 2015). High error-rate is also a likely explanation to the observed lower overall similarity of

343 locally collinear blocks (LCBs) in the miniasm part of our study (Figs. 4 and S2). Nevertheless, with ever-

344     improving sequencing chemistries and throughput the issue of high error-rates is likely to grow smaller in

345     the future. Error profile aware quality control methods could also help to counter the various miscalls

346     made by NGS platforms.

347     The actual sequencing coverage of PacBio reads has an apparently significant effect in the finalization of

348     the genomes: In the SSPACE-LongRead bacterial genome study (Boetzer & Pirovano, 2014), it was found

349     that long-read coverage from around 60X upwards did not further improve genome closure on the contig

350     level. Regarding read error-rates, it is already possible to self-correct PacBio CLR reads by using the reads

351     of insert pipeline of the SMRT Analysis toolkit. For each sequenced molecule, an improved consensus

352     sequence is obtained by aligning all the produced subreads together which cancels out the random errors

353     in individual reads. The final quality of the sequence depends on the number of subreads obtained for each

354     single molecule. Thanks to the nearly random error profile of the PacBio RS II instrument, single

355     nucleotide miscalls in the reads will not be propagated to the reads of insert output, that is, the circular

356     consensus (CCS) reads. Furthermore, the new Sequel instrument of Pacific Biosciences has promised 7-

357     fold throughput as compared to the earlier RS II platform, which has major ramifications also for the

358     throughput of corrected reads from the platform.

359     There may be additional approaches to the gap filling problem. In theory, a simple gap-tolerant alignment

360     of sequencing reads of variable lengths using existing mapping tools would be able to reliably span at

361     least short gaps, say 1-20 bp in length. This is one of the intriguing prospects of solving the gap filling

362     problem, especially as the average read lengths of next-generation sequencing platforms are likely to only

363     increase in the future.

364     CONCLUSIONS

365 Here, we presented an automated pipeline to solve the gap filling problem using a combination of

366 SSPACE-LongRead (Boetzer & Pirovano, 2014) and FGAP (Piro *et al.*, 2014) in a controlled manner and

367 wrapping these methods together in a pipeline called gapFinisher. gapFinisher ensures the uniqueness of

368 the BLAST alignments returned by the FGAP algorithm by iterating through the read data one read and

369 one superscaffold at a time. As evident from the result statistics (Table 2) and the visualizations of the

370 draft genomes (Supplementary figures S1 and S2), the control provided by gapFinisher leads into efficient

371 but reliable gap filling. The effects appear to scale up in a large genome *de novo* assembly (Table 2 and

372 Fig. 5).

373 The applicability of gapFinisher is currently limited to SSPACE-LongRead academic license version

374 output only and requires the user to be able to run SSPACE-LongRead at least once. Nevertheless,

375 SSPACE-LongRead currently is the only publicly available scaffolding software that is able to produce

376 information about the sequences spanning the gaps in the final scaffolds, i.e. the "inner-scaffold

377 sequences" subdirectory. Should other utilities with this key feature become available, we will further

378 develop gapFinisher for full compatibility. Our pipeline contributes to filling long gaps and solving the

379 non-trivial task of gap filling after scaffolding draft genomes of multiple organisms. Applying gapFinisher

380 will accelerate the finishing of draft genomes of both prokaryote and eukaryote organisms, even in

381 published genome assemblies.

382 The script to run the gapFinisher pipeline is published under GNU's general public license and can be

383 downloaded at: http://www.helsinki.fi/~jkammone/gapFinisher.zip

384 Competing interests: None declared.

385 ACKNOWLEDGEMENTS

386  The authors would like to acknowledge all supporting organizations and collaborators. The authors would

387  like to thank Dr. Ari Löytynoja for constructive comments on the manuscript text.

388  REFERENCES

389  Altschul SF, Gish W, Miller W *et al.* Basic local alignment search tool. *Journal of Molecular Biology*,

390  1990;215(3):403-10.

391  Bankevich A, Nurk S, Antipov D *et al.* SPAdes: a new genome assembly algorithm and its applications to

392  single-cell sequencing. *Journal of Computational Biology*, 2012;19(5): 455-477.

393  Boetzer M, Henkel CV, Jansen HJ *et al.* Scaffolding pre-assembled contigs using SSPACE.

394  *Bioinformatics* 2011;4(27): 578-579.

395  Boetzer M & Pirovano W. Toward almost finished genomes with GapFiller. *Genome Biology* 2012;13(6):

396  R56.

397  Boetzer M & Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes using long read

398  sequence information. *BMC Bioinformatics* 2014;15(1): 211.

399  Camacho C, Madden T, Ma N *et al.* BLAST command line applications user manual. National Center for

400  Biotechnology Information (US). http://www.ncbi.nlm.nih.gov/books/NBK1763/. (09 October 2017, date

401  last accessed)

402  Chaisson MJ & Tessler G. Mapping single molecule sequencing reads using basic local alignment with

403  successive refinement (BLASR): application and theory. *BMC Bioinformatics* 2012;13:238.

404  Christiansen T, Orwant J, Wall L, Foy B. Programming Perl. O'Reilly Media 2012.

405  Darling ACE, Mau B, Blattner FR & Perna NT. Mauve: Multiple Alignment of Conserved Genomic

406  Sequence With Rearrangements. *Genome Research*, 2004;14(7): 1394–1403.

407  Dohm JC, Lottaz C, Borodina T & Himmelbauer H. Substantial biases in ultra-short read data sets from

408  high-throughput DNA sequencing. *Nucleic Acids Research* 2008;16(36): e105.

409  Eidam C, Poehlein A, Brenner Michael G *et al.* Complete Genome Sequence of *Mannheimia haemolytica*

410  Strain 42548 from a Case of Bovine Respiratory Disease. *Genome announcements* 2013;1(3): e00318-13.

411  English AC, Richards S, Han Y *et al*. Mind the gap: upgrading genomes with Pacific Biosciences RS long-

412  read sequencing technology. *PloS ONE*, 2012;7(11), e47768.

413   Gentzsch W. Sun Grid Engine: Towards Creating a Compute Power Grid. In: *CCGRID '01: Proceedings*
414   *of the 1st International Symposium on Cluster Computing and the Grid.* 2001;35.

415   Harhay GP, McVey DS, Koren S *et al.* Complete Closed Genome Sequences of Three *Bibersteinia*
416   *trehalosi* Nasopharyngeal Isolates from Cattle with Shipping Fever. *Genome announcements* 2014;2(1):
417   e00084-14.

418   Kammonen JI, Smolander OP, Sipilä T *et al.* Increased transcriptome sequencing efficiency with modified
419   Mint-2 digestion-ligation protocol. *Analytical Biochemistry*, 2015;477:38-40.

420   de Koning AJ, Gu W, Castoe TA *et al.* Repetitive elements may comprise over two-thirds of the human
421   genome. *PLoS Genetics*, 2011;7(12), e1002384.

422   Langmead B & Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods,* 2012;9(4):357-
423   359.

424   Laver T, Harrison J, O'Neill PA *et al.* Assessing the performance of the Oxford Nanopore Technologies
425   MinION. *Biomolecular Detection and Quantification* 2015;3(3):1-8.

426   Li H. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences.
427   *Bioinformatics* 2016;32(14):2103-10.

428   Li H & Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
429   *Bioinformatics*, 2009;25(14):1754-1760.

430   Li YI & Copley RR. Scaffolding low quality genomes using orthologous protein sequences.
431   *Bioinformatics* 2013;29(2): 160-165.

432   MATLAB and Statistics Toolbox Release 2015a, The MathWorks, Inc., Natick, Massachusetts, United
433   States.

434   Magoč T & Salzberg SL. FLASH: fast length adjustment of short reads. *Bioinformatics* 2011;27(21):
435   2957-2963.

436   Margulies M, Egholm M, Altman WE *et al*. Genome sequencing in microfabricated high-density picolitre
437   reactors. *Nature,* 2005;437:376–380.

438   Nakamura K, Oshima T, Morimoto T *et al.* Sequence-specific error profile of Illumina sequencers.
439   *Nucleic Acids Research*, 2011;13(39): e90.

440   Piro VC, Faoro H, Weiss VA *et al*. FGAP: an automated gap closing tool. *BMC Research Notes*
441   2014;7:371.

442     Salmela L, Sahlin K, Mäkinen V & Tomescu A. Gap Filling as Exact Path Length Problem. *Journal of*
443     *Computational Biology* 2016;23(5):347-61.

444     Schirmer M, Ijaz UZ, D'Amore R *et al.* Insight into biases and sequencing errors for amplicon sequencing
445     with the Illumina MiSeq platform. *Nucleic Acids Research*, 2015;6(43), e37.

446     Slater GS & Birney E. Automated generation of heuristics for biological sequence comparison. *BMC*
447     *Bioinformatics* 2005;6:31.

448     Smit AFA, Hubley R & Green P. 2013-2015. *RepeatMasker Open-4.0*. Retrieved from: Smit, AFA,
449     Hubley, R & Green, P. RepeatMasker Open-4.0.
450     2013-2015 <http://www.repeatmasker.org>http://www.repeatmasker.org (15 November 2017, date last
451     accessed)

452     Walker BJ, Abeel T, Shea T *et al.* Pilon: An Integrated Tool for Comprehensive Microbial Variant
453     Detection and Genome Assembly Improvement. *PLoS ONE* 2014;9(11): e112963.

454    FIGURE CAPTIONS

455    **Figure 1.** *Visualization of the FGAP* (Piro *et al.*, 2014) *and gapFinisher workflows.* **a)** FGAP is expected

456    to find identity between a long read (blue bar) and two contigs (gray blocks) separated by a gap (N's) and

457    then fill the gaps with the sequence. **b)** In practice, FGAP is allowed to find multiple places for one long

458    read and nonhomologous gaps with the same sequence. **c)** gapFinisher uses the association of the long

459    read and the scaffold reported by BLASR (Chaisson & Tessler, 2012) of SSPACE-LongRead (Boetzer &

460    Pirovano, 2014) and ensures that each long read is only used once in gap filling.

461    **Figure 2.** *A more detailed visualization of the gapFinisher pipeline workflow.* **a)** SSPACE-LR (Boetzer &

462    Pirovano, 2014) reports new scaffolds and these are iterated through one scaffold at a time. **b)** SSPACE-

463    LR output shows the PacBio reads associated with the gaps in the scaffolds. **c)** These reads are then

464    circulated through the FGAP (Piro *et al*., 2014) pipeline with only the single scaffold as input data. This

465    logical step prevents same PacBio reads from being used in parts of the draft genome other than the

466    current scaffold. Measures are then taken to either **d1)** replace the unknown sequence with that of the long

467    read **(**=fill gap) or **d2)** reject the alignment and leave the gap to the genome as is.

468    **Figure 3.** *minidot* (Li, 2016) *plots of the Mannheimia haemolytica draft genome at different stages of the*

469    *assembly.* **Top top left:** Image key and reading direction. **Top left:** SPAdes contig assembly, **top center:**

470    scaffold stage of SPAdes contigs, **top right:** gap filling stage, of the *M. haemolytica* draft genome.

471    **Bottom left:** miniasm (Li, 2016) contig assembly, **bottom center:** scaffold stage of miniasm contigs,

472    **bottom right:** gap filling stage, of the *M. haemolytica* draft genome. The red diagonal lines indicate

473    continuous regions of alignment between the draft assembly and the *M. haemolytica* reference sequence.

474    The blue diagonal lines indicate regions with inverted alignment. The red and blue dots indicate repeats

475    and inverted repeats, respectively. Draft assembly contig/scaffold boundaries are shown as grey vertical

476    lines. The alignment plots are provided for each of the bacterial genomes as supplementary figure S1.

477 **Figure 4.** *Mauve* (Darling *et al.*, 2004) *alignments of the Mannheimia haemolytica genome*. The

478 visualizations are from **a)** before and **b)** after the scaffolding/gap filling stage. The corresponding Locally

479 Collinear Blocks (LCB) in the three genome versions are indicated by different colors of horizontal bars.

480 The darker lines within the blocks indicate local median similarity while the light lines show the range of

481 local similarity values. White areas indicate low or no similarity. Blocks below the center line indicate

482 regions that align in the reverse complement (inverse) orientation. **a):** *M. haemolytica* reference sequence

483 (**red bar**), SPAdes (Bankevich *et al.*, 2012) assembly contig sequences (**green bar**), and miniasm (Li,

484 2016) assembly contig sequences (**blue bar**). **b):** *M. haemolytica* reference sequence (**red bar**), and gap

485 filled scaffolds using the SPAdes assembly contig sequences (**green bar**), and the miniasm assembly

486 contig sequences (**blue bar**).

487 **Figure 5.** *Results from the performance benchmark of the assembly, scaffolding and gap filling tools used.*

488 The exact values are reported in supplementary table S3. **a)** Peak random access memory (RAM) use in

489 gigabytes (GB) in the six bacterial assemblies. **b)** Runtimes (in minutes) of the bacterial assemblies. The

490 *F. tularensis* and *S. enterica* assembly runtimes are omitted due to the large number of CLR reads per

491 scaffold reported in the gap filling stage (953 and 317, respectively) and the consequent long runtimes. **c)**

492 Peak RAM use (GB) and runtimes (in hours) of the assembly, scaffolding and gap filling algorithms in the

493 marine mammal (*Phocidae* family) genome assembly.

494 **Figure S1.** *minidot* (Li, 2016) *plots of the six bacterial genomes at different stages of the assembly*. **a)** *E.*

495 *coli* K12, **b)** *E. coli* O157:H7, **c)** *B. trehalosi*, **d)** *M. haemolytica*, **e)** *F. tularensis*, **f)** *S. enterica*. **Top top**

496 **left:** Image key and reading direction. **Top row** (in all subfigures): SPAdes (Bankevich *et al.*, 2012) contig

497 assembly, scaffolding and gap filling (gapFinisher) stages of the assembly. **Bottom row** (in all subfigures):

498    miniasm (Li, 2016) contig assembly, scaffolding and gap filling (gapFinisher) stages of the assembly. The

499    scaffolding and gap filling stages are missing for the *E. coli* K12 assembly (**a)**) since the genome was in a

500    single sequence (i.e. closed) after miniasm.

501    **Figure S2.** *Mauve* (Darling *et al*., 2004) *alignments of the six bacterial genomes at different stages of the*

502    *assembly.* **a)** *E. coli* K12, **b)** *E. coli* O157:H7, **c)** *B. trehalosi*, **d)** *M. haemolytica*, **e)** *F. tularensis*, **f)** *S.*

503    *enterica*. **Top part** (in all subfigures): progressiveMauve alignment of the respective bacterial reference

504    genome (**red bar**), the SPAdes (Bankevich *et al*., 2012) contig draft genome (**green bar**) and the miniasm

505    (Li, 2016) contig draft genome (**blue bar**). **Bottom part** (in all subfigures): progressiveMauve alignment

506    of the respective bacterial reference genome (**red bar**), the SPAdes assembly pipeline gap filled

507    (gapFinisher) scaffolds (**green bar**) and the miniasm assembly pipeline gap filled (gapFinisher) scaffolds

508    (**blue bar**). Only the contig assembly stage (top part) is shown for the *E. coli* K12 assembly (subfigure **a)**)

509    since the genome was in a single sequence (i.e. closed) after miniasm.

510 TABLE CAPTIONS

511 **Table 1.** *Next-generation sequencing read statistics and sequencing coverage for the sample datasets*. The

512 bacterial data are from 2013 and originate from the Sequencing Read Archive (SRA). The sequencing

513 chemistries were not accurately described in the original datasets but the bacterial MiSeq read data

514 represent either Illumina (San Diego, California, United States) sequencing-by-synthesis chemistry v1 or

515 v2. The mammal MiSeq read data are a mixture of Illumina sequencing-by-synthesis chemistry v2 and v3.

516 The bacterial Pacific Biosciences (Menlo Park, California, United States) RS reads represent PacBio

517 SMRT sequencing chemistries that are earlier than P4-C2 and the mammal PacBio RS reads are a mixture

518 of PacBio SMRT sequencing chemistries P5-C3 and P6-C4.


519 **Table 2.** *De novo assembly, scaffolding and gap filling statistics for the model genomes.* For clarity, only

520 the most significant statistics are shown here and the full statistics provided as supplementary table S3.


521 **Table S3.** *All de novo assembly, scaffolding and gap filling statistics for the model genomes*. In addition,

522 the performance benchmark statistics are included in the last three columns.

523    **Table S4.** *Gap filling data used and FGAP* (Piro *et al*., 2014) *default test results reported for an*

524    *unpublished draft genome of a marine mammal from the Phocidae* f*amily*. An admittedly small number of

525    Pacific Biosciences RS II platform circular consensus reads (*reads of insert*) with summed length of about

526    280 kbp filled 45.5 million unknown bases in the draft genome, a result reported by FGAP with the default

527    alignment settings. By further changing the FGAP command line options, one can adjust the number of

528    BLAST (Altschul *et al*., 1990) hits returned to perform the alignment. By default, this is 200 hits. We

529    further ran another test, where we reduced this amount to two so that only the best two BLAST hits would

530    be considered in the gap filling. Still, more than 4.5 million N's were reportedly filled with our test set, a

531    far greater number of bases than contained by the original read data used. The default BLAST alignment

532    parameters of FGAP for opening and extending a gap are both set to the value 1. The default values in

533    command line applications of BLAST for opening and extending a gap are set as 5 and 2, respectively, as

534    written in the BLAST user guide by Camacho *et al*. Thus, depending on the total score of the alignment,

535    gap opening in the alignment is up to 80% and gap extension up to 50% more likely than the BLAST

536    defaults. The minimum raw score of a BLAST hit in FGAP is set to value 25, a typical raw score value of

537    highly dissimilar sequences irrespective of the gap penalty parameters. Moreover, a maximum of 200

538    BLAST results may be returned for a 70 percent identity in alignment length of 300 bp by default. In

539    general, the default parameters of FGAP appear too lenient and may fill gaps based on alignments that are

540    incorrect and may appear multiple times where unique alignments are desired.

# Figure 1

Visualization of the FGAP (Piro *et al.*, 2014) and gapFinisher workflows.

**a)** FGAP is expected to find identity between a long read (blue bar) and two contigs (gray blocks) separated by a gap (N's) and then fill the gaps with the sequence. **b)** In practice, FGAP is allowed to find multiple places for one long 5 read and nonhomologous gaps with the same sequence. **c)** gapFinisher uses the association of the long read and the scaffold reported by BLASR (Chaisson & Tessler, 2012) of SSPACE-LongRead (Boetzer & Pirovano, 2014) and ensures that each long read is only used once in gap filling.
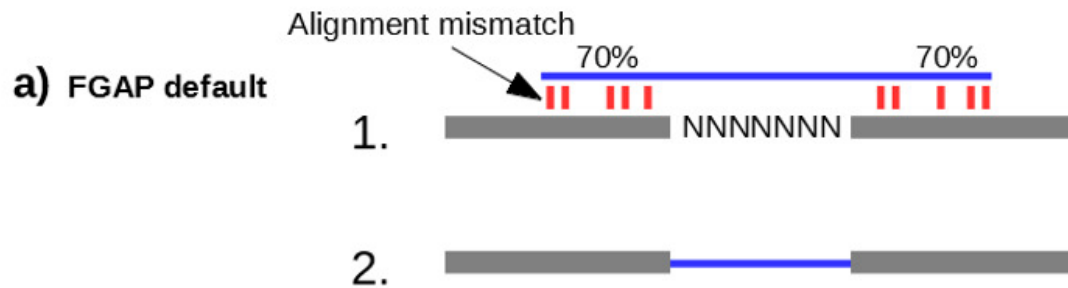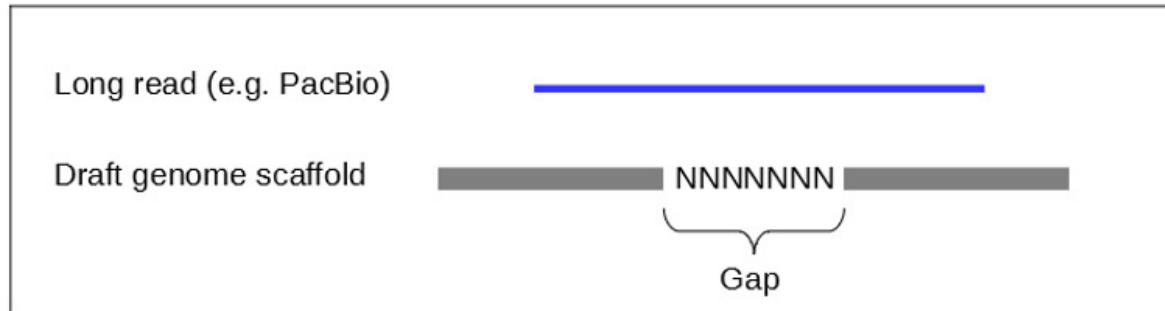
# Figure 2

A more detailed visualization of the gapFinisher pipeline workflow.

**a)** SSPACE-LR (Boetzer & Pirovano, 2014) reports new scaffolds and these are iterated through one scaffold at a time. **b)** SSPACE-LR output shows the long reads associated with the gaps in the scaffolds. **c)** These reads are then circulated through the FGAP (Piro *et al.*, 2014) pipeline with only the single scaffold as input data. This logical step prevents same reads from being used in parts of the draft genome other than the current scaffold. Measures are then taken to either **d1)** replace the unknown sequence with that of the long read (=*fill gap*) or **d2)** reject the alignment and leave the gap to the genome as is.
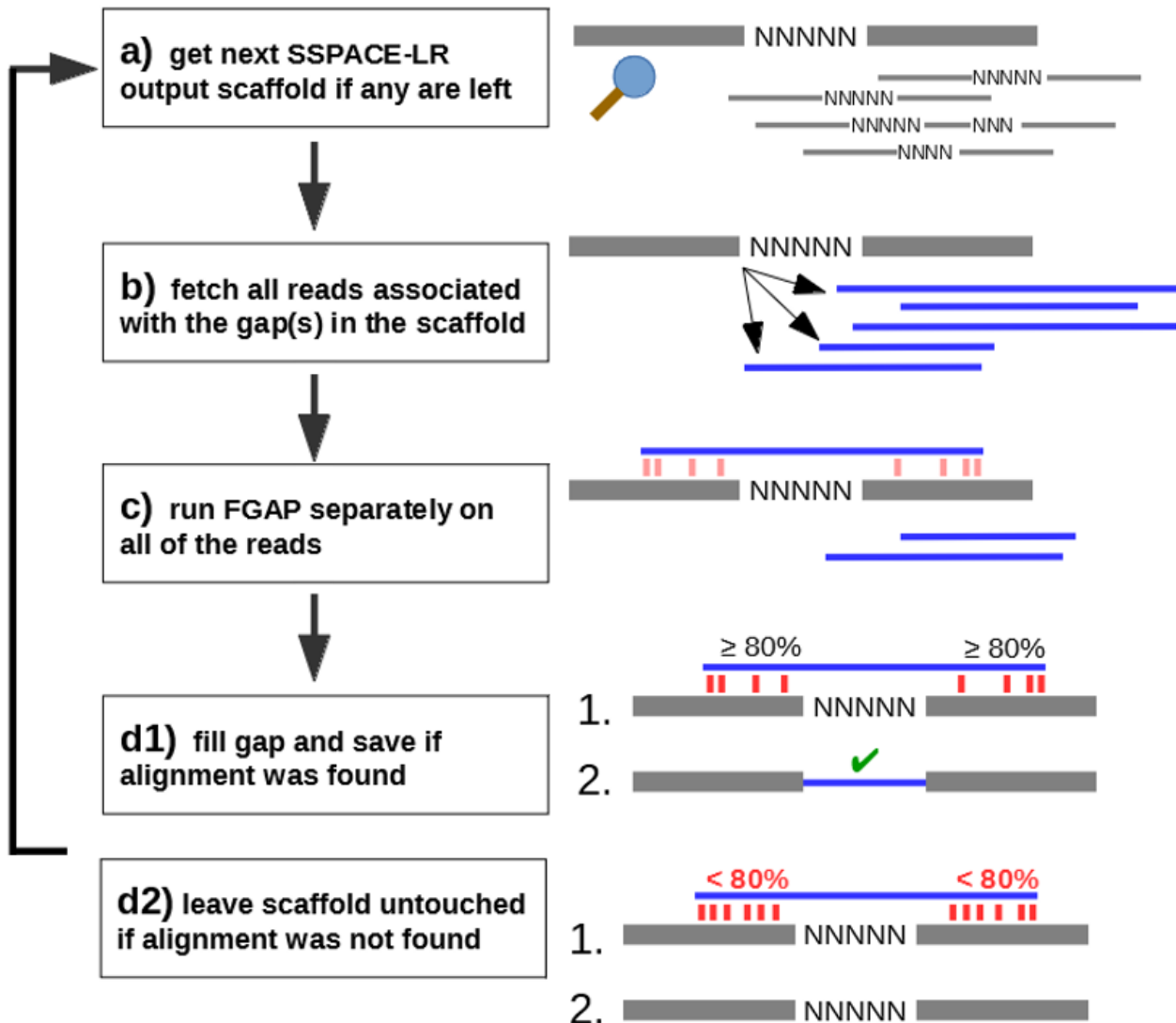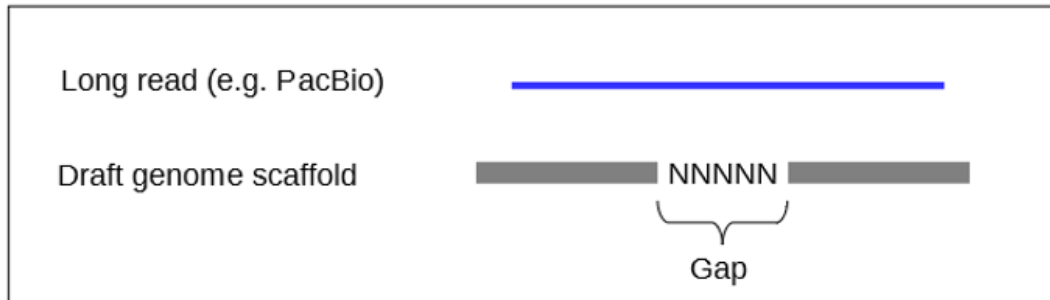
# Figure 3

minidot (Li, 2016) plots of the *Mannheimia haemolytica* draft genome at different stages of the assembly.
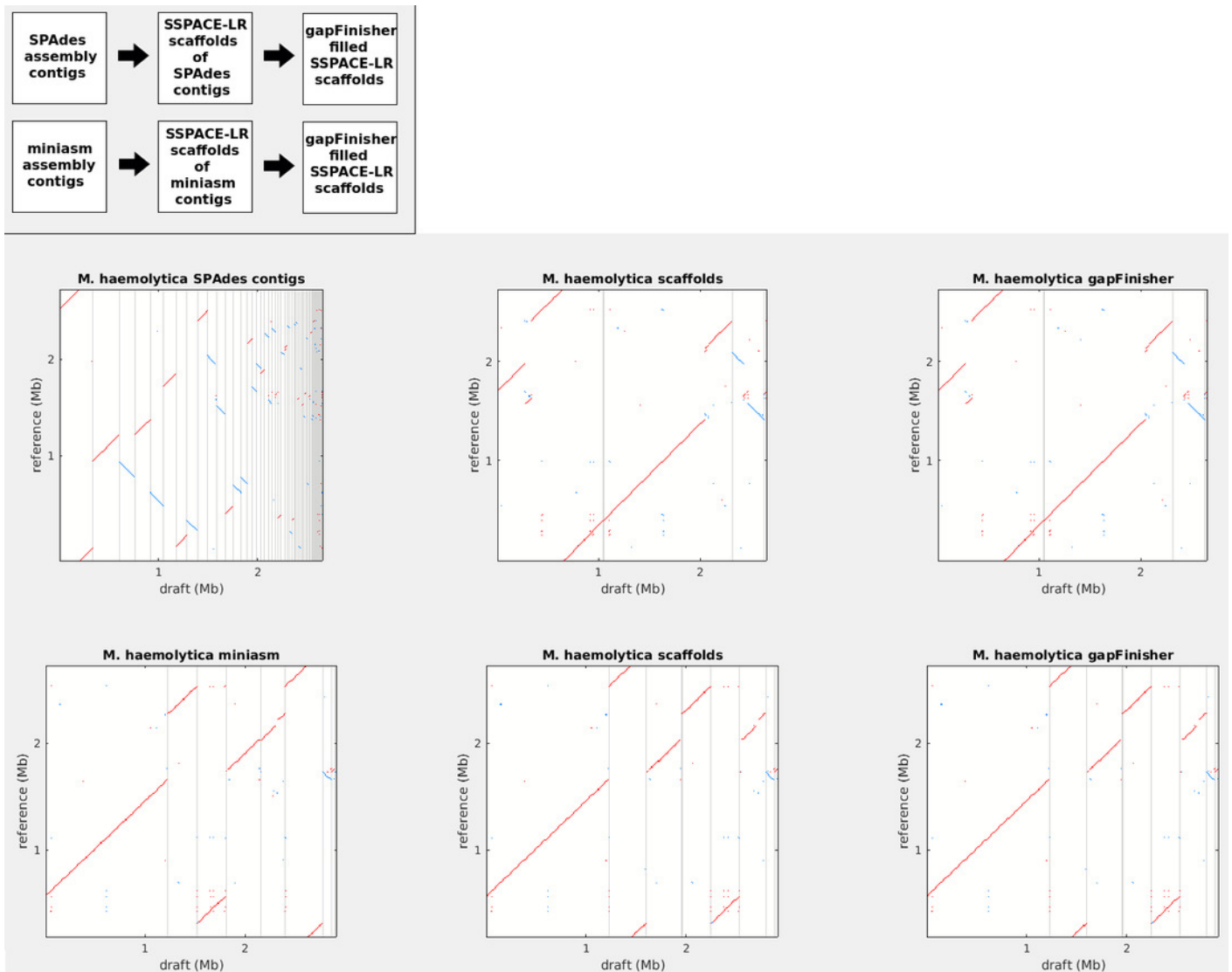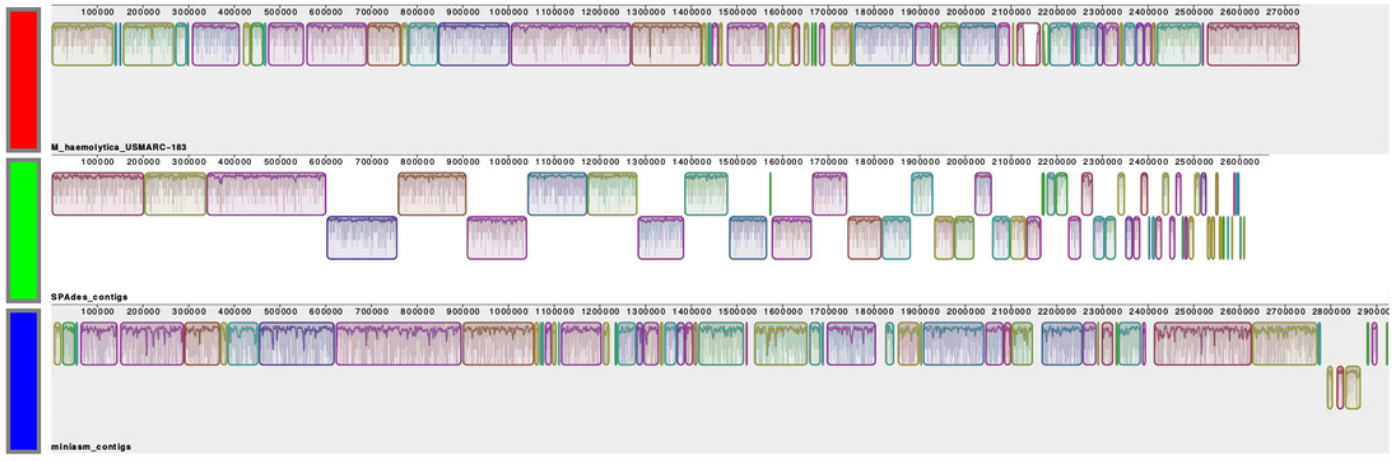
# Figure 4

Mauve (Darling *et al.*, 2004) alignments of the *Mannheimia haemolytica* genome.

The visualizations are from **a)** before and **b)** after the scaffolding/gap filling stage. The corresponding Locally Collinear Blocks (LCB) in the three genome versions are indicated by different colors of horizontal bars. The darker lines within the blocks indicate local median similarity while the light lines show the range of local similarity values. White areas indicate low or no similarity. Blocks below the center line indicate regions that align in the reverse complement (inverse) orientation. **a)** *M. haemolytica* reference sequence (red bar), SPAdes (Bankevich *et al.*, 2012) assembly contig sequences (green bar), and miniasm (Li, 2016) assembly contig sequences (blue bar). **b)** *M. haemolytica* reference sequence (red bar), and gap filled scaffolds using the SPAdes assembly contig sequences (green bar), and the miniasm assembly contig sequences (blue bar).
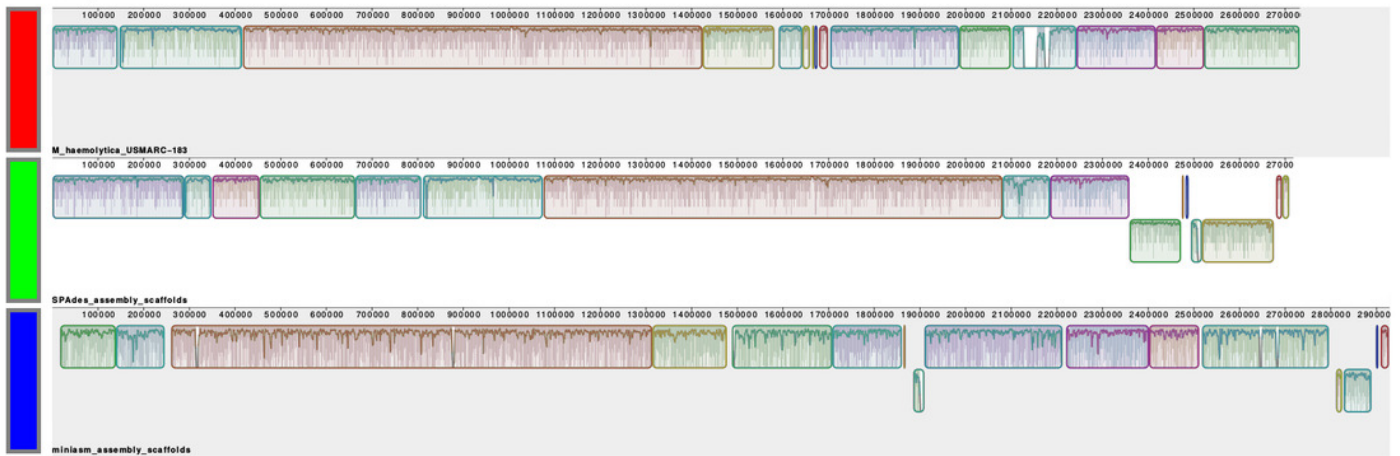
# Figure 5

Results from the performance benchmark of the assembly, scaffolding and gap filling tools used.

a) Peak RAM use in bacterial assemblies

b) Runtimes of bacterial assemblies

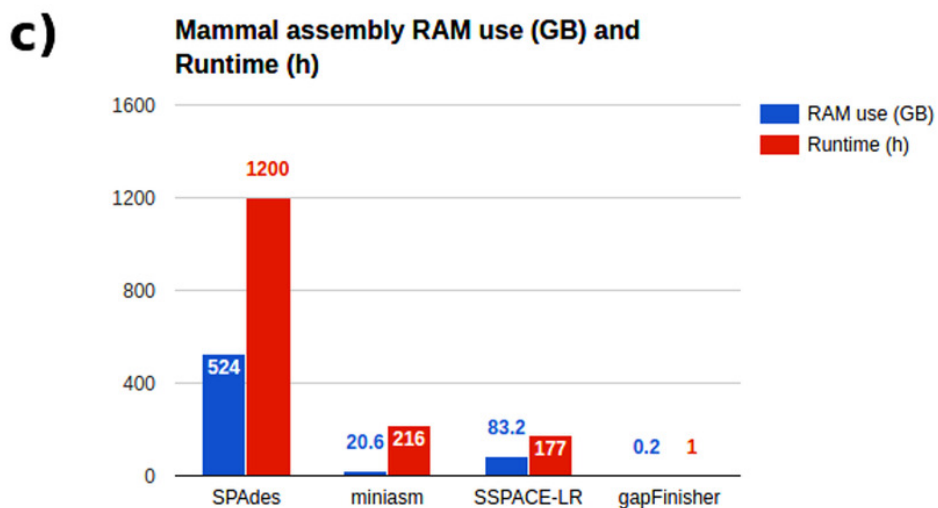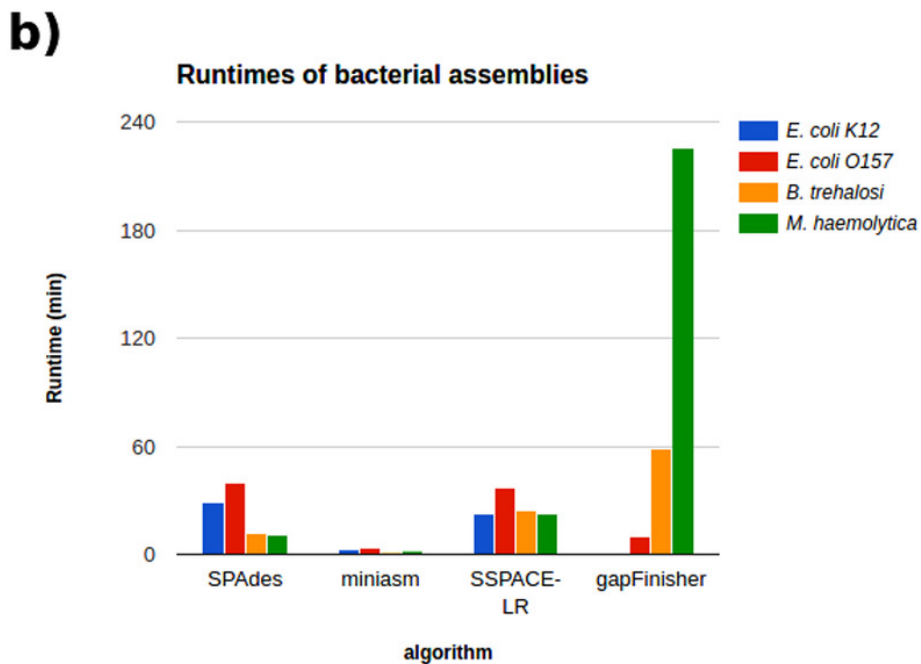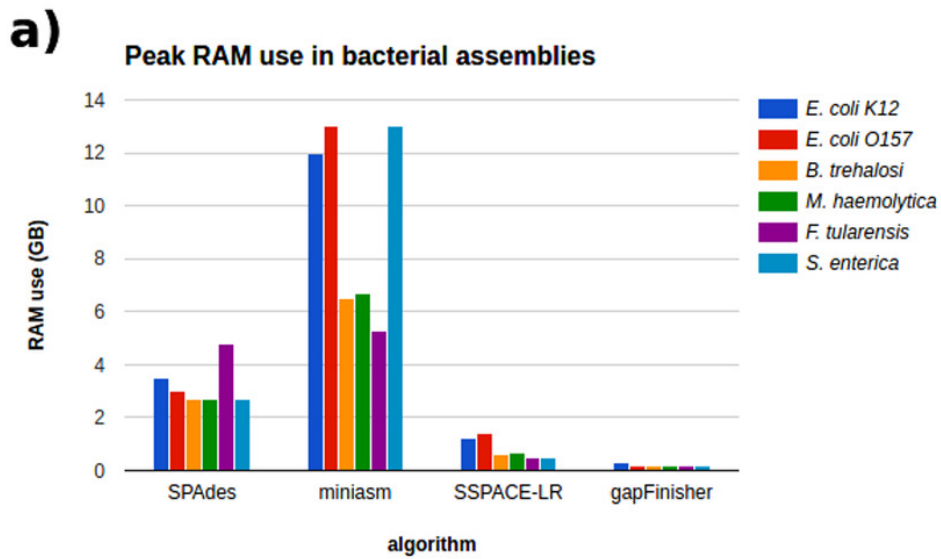c) Mammal assembly RAM use (GB) and Runtime (h)

**Peer**J Preprints

# Table 1 (on next page)

Next-generation sequencing read statistics and sequencing coverage for the sample datasets.

| | Illumina MiSeq paired end reads (100X) | | | PacBio RS reads (200X) | | |
|---|---|---|---|---|---|---|
| Organism | Total reads | Total bases | Avg. read length | Total reads | Total bases | Avg. read length |
| E. coli K12 MG1655 | 3,046,358 | 460,000,058 | 151 | 383,482 | 929,129,994 | 2,422 |
| E. coli O157:H7 | 3,794,862 | 548,505,079 | 144 | 403,919 | 1,100,295,861 | 2,724 |
| B. trehalosi | 1,718,212 | 249,216,010 | 145 | 205,096 | 499,939,066 | 2,437 |
| M. haemolytica | 1,724,414 | 249,368,724 | 144 | 175,953 | 531,234,319 | 3,019 |
| F. tularensis | 926,716 | 199,169,591 | 214 | 176,376 | 399,767,452 | 2,266 |
| S. enterica | 1,943,848 | 279,774,061 | 143 | 394,699 | 1,000,244,555 | 2,534 |
| Organism | Illumina MiSeq reads (25X) | | | PacBio RS reads (50X) | | |
| Mammal | 329,484,322 | 62,120,890,467 | 188 | 17,695,174 | 146,961,409,902 | 8,305 |

# Table 2(on next page)

*De novo* assembly, scaffolding and gap filling statistics for the model genomes.

| | | Sequences | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Organism | Tool | Expected | Final | Sum (bp) | N50 length | Gaps (bp) | Gap % | LCBs* | Indels | Inversions |
| *E. coli K12* | SPAdes | 1 | 35 | 4,661,027 | 4,640,853 | 0 | 0.00% | 4 | 0 | 2 |
| *MG1655* | SSPACE-LR | 1 | 34 | 4,661,028 | 4,641,005 | 1 | 0.00% | 4 | 0 | 2 |
| | gapFinisher | | | 4,661,028 | 4,641,005 | 1 | 0.00% | | | |
| | miniasm | 1 | **1** | 4,793,967 | 4,793,967 | 0 | 0.00% | 4 | 0 | 2 |
| | SSPACE-LR | - | - | - | - | - | - | - | - | - |
| | gapFinisher | | | - | - | - | | | | |
| *E. coli* | SPAdes | 10 | 87 | 5,547,646 | 3,323,349 | 3 | 0.00% | 27 | 5 | 5 |
| *O157:H7* | SSPACE-LR | 10 | 50 | 5,568,199 | 3,323,349 | 18486 | 0.33% | 23 | 4 | 3 |
| | gapFinisher | | | 5,568,974 | 3,323,349 | **5750** | 0.10% | | | |
| | miniasm | 10 | 25 | 5,898,494 | 537,223 | 0 | 0.00% | 27 | 7 | 7 |
| | SSPACE-LR | 10 | 16 | 5,908,008 | 61,209 | 9514 | 0.16% | 23 | 6 | 5 |
| | gapFinisher | | | 5,907,537 | 61,209 | **2495** | 0.04% | | | |
| *B. trehalosi* | SPAdes | 1 | 51 | 2,376,258 | 274,711 | 2 | 0.00% | 30 | 7 | 8 |
| | SSPACE-LR | 1 | 12 | 2,401,287 | 438,635 | 2804 | 0.12% | 21 | 6 | 3 |
| | gapFinisher | | | 2,401,265 | 438,599 | **401** | 0.02% | | | |
| | miniasm | 1 | 17 | 2,510,680 | 221,473 | 0 | 0.00% | 30 | 3 | 7 |
| | SSPACE-LR | 1 | 10 | 2,520,563 | 377,524 | 9883 | 0.39% | 21 | 4 | 5 |
| | gapFinisher | | | 2,521,341 | 37,752 | **4920** | 0.20% | | | |
| *M. haemolytica* | SPAdes | 1 | 112 | 2,664,209 | 101,958 | 35 | 0.00% | 40 | 5 | 7 |
| | SSPACE-LR | 1 | 17 | 2,718,326 | 1,073,880 | 13504 | 0.50% | 13 | 4 | 2 |
| | gapFinisher | | | 2,717,906 | 1,073,740 | **4498** | 0.17% | | | |
| | miniasm | 1 | 10 | 2,926,783 | 378,549 | 0 | 0.00% | 40 | 4 | 1 |
| | SSPACE-LR | 1 | 8 | 2,928,560 | 378,549 | 1777 | 0.06% | 13 | 4 | 1 |
| | gapFinisher | | | 2,928,834 | 378,549 | **1155** | 0.04% | | | |
| *F. tularensis* | SPAdes | 3 | 135 | 1,807,729 | 25,688 | 0 | 0.00% | 80 | 12 | 18 |
| | SSPACE-LR | 3 | 58 | 1,855,045 | 56,838 | 23176 | 1.25% | 42 | 9 | 9 |
| | gapFinisher | | | 1,851,864 | 56,791 | **11254** | 0.61% | | | |
| | miniasm | 3 | 20 | 2,000,228 | 15,305 | 0 | 0.00% | 80 | 5 | 4 |
| | SSPACE-LR | 3 | 9 | 2,021,978 | 426,098 | 21750 | 1.08% | 42 | 4 | 4 |
| | gapFinisher | | | 2,021,618 | 425,969 | **16843** | 0.83% | | | |
| *S. enterica* | SPAdes | 4 | 217 | 4,982,997 | 153,597 | 655 | 0.01% | 49 | 4 | 10 |
| | SSPACE-LR | 4 | 94 | 5,026,381 | 1,020,795 | 10050 | 0.20% | 27 | 3 | 6 |
| | gapFinisher | | | 5,028,882 | 1,020,937 | **2917** | 0.06% | | | |
| | miniasm | 4 | 16 | 5,373,212 | 735,723 | 0 | 0.00% | 49 | 5 | 2 |
| | SSPACE-LR | 4 | 10 | 5,384,667 | 874,322 | 11455 | 0.21% | 27 | 5 | 2 |
| | gapFinisher | | | 5,384,057 | 874,322 | **4641** | 0.09% | | | |
| *Mammal* | SPAdes | Unkn. | 236592 | 2,253,617,865 | 19,739 | 0 | 0.00% | - | - | - |
| | SSPACE-LR | Unkn. | 10143 | 2,462,623,627 | 599,108 | 10136364 | 0.41% | - | - | - |
| | gapFinisher | | | 2,466,785,189 | 601,444 | **6945295** | 0.28% | | | |
| | miniasm | Unkn. | 1314 | 2,460,097,408 | 8,668,858 | 0 | 0.00% | - | - | - |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| S-LR / 1st | Unkn. | 1115 | 2,460,626,045 | 9,381,548 | 528637 | 0.02% | - | - | - |
| gapF / 1st | | | 2,460,674,964 | 9,381,548 | **351878** | 0.01% | | | |
| S-LR / 2nd | Unkn. | 1008 | 2,460,965,525 | 9,562,827 | 642439 | 0.03% | - | - | - |
| gapF /2nd | | | 2,460,993,827 | 9,562,827 | **616617** | 0.03% | | | |

*LCB = Locally Collinear Block of progressiveMauve (27) alignment.