

1 **Title**

2 SNVPhyl: A Single Nucleotide Variant Phylogenomics pipeline for microbial genomic epidemiology

3 **Authors**

4 Aaron Petkau^{1,*}, Philip Mabon¹, Cameron Sieffert¹, Natalie Knox¹, Jennifer Cabral¹, Mariam Iskander^{1,2},
5 Mark Iskander^{1,2}, Kelly Weedmark^{1,3}, Rahat Zaheer^{1,4}, Lee S. Katz⁵, Celine Nadon¹, Aleisha Reimer¹,
6 Eduardo Taboada⁶, Robert G. Beiko⁷, William Hsiao^{8,9}, Fiona Brinkman¹⁰, Morag Graham^{1,2}, the IRIDA
7 Consortium, and Gary Van Domselaar^{1,2}

8 ¹National Microbiology Laboratory, Public Health Agency of Canada, Winnipeg, Manitoba R3E 3R2,
9 Canada

10 ²University of Manitoba, Winnipeg, Manitoba R3T 2N2, Canada

11 ³Health Canada - Bureau of Microbial Hazards, Ottawa, Ontario K1A 0K9, Canada

12 ⁴Lethbridge Research and Development Centre, Lethbridge, Alberta T1J 4B1, Canada

13 ⁵Centers for Disease Control and Prevention, Atlanta, Georgia 30333, USA

14 ⁶National Microbiology Laboratory, Public Health Agency of Canada, Lethbridge, Alberta T1J 3Z4,
15 Canada

16 ⁷Dalhousie University, Halifax, Nova Scotia B3H 4R2, Canada

17 ⁸BC Public Health Microbiology and Reference Laboratory, Vancouver, British Columbia V5Z 4R4,
18 Canada

19 ⁹University of British Columbia, Vancouver, British Columbia V6T 1Z4, Canada

20 ¹⁰Simon Fraser University, Burnaby, British Columbia V5A 1S6, Canada

21 *Corresponding author (aaron.petkau@phac-aspc.gc.ca)

22 **Abstract**

23 **Motivation:** The recent widespread application of whole-genome sequencing (WGS) for microbial
24 disease investigations has spurred the development of new bioinformatics tools, including a notable
25 proliferation of phylogenomics pipelines designed for infectious disease surveillance and outbreak
26 investigation. Transitioning the use of WGS data out of the research lab and into the front lines of
27 surveillance and outbreak response requires user-friendly, reproducible, and scalable pipelines that have
28 been well validated.

29 **Results:** SNVPhyl (Single Nucleotide Variant Phylogenomics) is a bioinformatics pipeline for identifying
30 high-quality SNVs and constructing a whole genome phylogeny from a collection of WGS reads and a
31 reference genome. Individual pipeline components are integrated into the Galaxy bioinformatics
32 framework, enabling data analysis in a user-friendly, reproducible, and scalable environment. We show
33 that SNVPhyl can detect SNVs with high sensitivity and specificity and identify and remove regions of
34 high SNV density (indicative of recombination). SNVPhyl is able to correctly distinguish outbreak from
35 non-outbreak isolates across a range of variant-calling settings, sequencing-coverage thresholds, or in the
36 presence of contamination.

37 **Availability:** SNVPhyl is available as a Galaxy workflow, Docker and virtual machine images, and a
38 Unix-based command-line application. SNVPhyl is released under the Apache 2.0 license and available at
39 <http://snvphyl.readthedocs.io/> or at <https://github.com/phac-nml/snvphyl-galaxy>.

40 **Introduction**

41 The high-efficiency and cost-effectiveness of whole-genome sequencing (WGS) using next-generation
42 sequencing (NGS) technologies is transforming the biomedical landscape. Entire microbial genomes can
43 be rapidly sequenced and subsequently queried with nucleotide-level resolution, an exciting new ability

44 that far outstrips other traditional microbial typing methods. This powerful new ability has the potential to
45 advance many fields, including in particular the field of infectious disease genomic epidemiology. A
46 number of landmark studies have demonstrated the power of WGS for molecular epidemiology. One
47 notable study is the investigation into the 2010 Haiti cholera outbreak (1-3), where WGS and
48 epidemiological data was used in support of the hypothesis that cholera was introduced to Haiti from UN
49 peacekeepers originally infected in Nepal. WGS has supported the investigation of outbreaks of
50 organisms as diverse as *Mycobacterium tuberculosis* (4, 5), *Escherichia coli* (6), and *Legionella*
51 *pneumophila* (7). These high-profile successes have motivated public health institutions and food
52 regulatory agencies to incorporate WGS into their routine microbial infectious disease surveillance and
53 outbreak investigation activities. The GenomeTrakr network used by the Centers for Disease Control
54 (CDC) and the Food and Drug Administration (FDA) agencies in the United States (8), PulseNet
55 International (<http://www.cdc.gov/pulsenet/next-generation.html>), Statens Serum Institut in Denmark (9),
56 and Public Health England (10) are leading the charge in this area, and have incorporated a variety of
57 analytical approaches to integrate WGS into their infectious disease surveillance activities. Two
58 approaches in particular have emerged as feasible methods for bacterial genomic epidemiology: gene-by-
59 gene methods, which extend the idea of multilocus sequence typing (MLST) to encompass a given
60 organism's entire genome (whole-genome MLST, wgMLST) or core genome (core genome MLST,
61 cgMLST) (11, 12); and single nucleotide variant (SNV)-based methods, which identify variants by
62 comparing a population of target genomes against a reference (13, 14). Gene-by gene methods are
63 promising as they are more amenable to assigning consistent sequence types using standardized MLST
64 schemas, but these schemas must be developed and validated for each organism. SNV-based methods are
65 popular as they do not require development of MLST schemas, but the variability in SNV-identification
66 methods and reference genome selection means they do not yet produce standard sequence types useful
67 for global communication of circulating infectious disease (12, 14). Where applicable, these two methods
68 are often combined (15).

69 A growing number of SNV-based pipelines have been developed (Table 1) and are distributed in
70 the form of web services (16), command-line software (17), or both (14). Web services provide a user-
71 friendly method of running large-scale analyses but require the uploading of sequence reads and rely on
72 third-party computing infrastructure, which may be inadequate for the analysis of typically large datasets
73 or due to data privacy concerns. Locally installed pipelines avoid the transfer of large datasets to third-
74 party websites, offer greater control over the execution environment for reproducibility, and allow for the
75 incorporation into pre-existing bioinformatics analysis environments. However, locally installed pipelines
76 may require considerable expertise to operate and can have substantial computing requirements.
77 Additionally, for many SNV-based pipelines, recombination detection and removal may require pre-
78 analysis to identify phage and genomic islands on the reference genome, or post-analysis with
79 computationally intensive recombination-detection software such as Gubbins (18) or ClonalFrameML
80 (19) to identify and mask possible recombinant regions. While a large choice of pipelines is available, a
81 systematic comparison of popular SNV pipelines has demonstrated that they generate highly concordant
82 phylogenetic trees but with variation in the particular SNVs identified (20). However, variation in the
83 installation procedures and execution environments of these pipelines proves challenging for integration
84 into a larger bioinformatics analysis system.

85 Galaxy (21) is a web-based biological data analysis platform that can be accessed through a
86 publicly available website, a locally installed instance linked to a high-performance compute cluster, or a
87 cloud-based environment. Galaxy provides a user-friendly web interface for the construction of data
88 analysis workflows using a mixture of built-in or community developed bioinformatics tools.
89 Additionally, Galaxy provides an API for automated workflow execution or other automations via
90 external software. These features have encouraged some software developers to integrate Galaxy within
91 larger data analysis systems. Examples of such analysis systems include IRIDA (<http://irida.ca>), the
92 Refinery Platform (<http://www.refinery-platform.org/>), and the Genomics Virtual Laboratory (22).

93 The SNVPhyl pipeline provides a reference-based SNV discovery and phylogenomic tree-
94 building pipeline along with ancillary tools integrated within the Galaxy framework. SNVPhyl can
95 quickly analyze many genomes, identify variants and generate a maximum likelihood phylogeny, an all-
96 against-all SNV distance matrix, as well as additional quality information to help guide interpretation of
97 the results. The pipeline has been under continuous development and refinement at Canada's National
98 Microbiology Laboratory since 2010; it is currently being used for outbreak investigations and will be
99 part of the validated suite of tools used by PulseNet Canada for routine foodborne disease surveillance
100 activities. Here, we describe the overall operation of SNVPhyl, survey its advanced features such as
101 repeat and recombination masking, and demonstrate its SNV-calling and phylogenomic tree building
102 accuracy using simulated and real-world datasets.

103 **Methods**

104 **SNVPhyl pipeline**

105 The SNVPhyl pipeline (Figure 1) consists of a set of pre-existing and custom-developed bioinformatics
106 tools for reference mapping, variant discovery, and phylogeny construction from identified SNVs. Each
107 stage of the pipeline is implemented as a separate Galaxy tool and the stages are joined together to
108 construct the SNVPhyl workflow. Distribution of the dependency tools for SNVPhyl is managed through
109 the Galaxy Toolshed (23). Scheduling of each tool is managed by Galaxy, which provides support for
110 execution on a single machine, high-performance computing environments utilizing most major
111 scheduling engines (e.g., Slurm, TORQUE, Open Grid Engine), or cloud-based environments.

112 **Input**

113 SNVPhyl requires as input a set of microbial WGS datasets, a reference genome, and an optional masking
114 file defining regions on the reference genome to exclude from the analysis. The sequencing data consists
115 of either single-end or paired-end reads. The reference genome consists of a draft or finished genome,
116 chosen typically to have high similarity with the collection of genome sequences under analysis. The

117 masking file stores the sequence identifier of the reference genome and the coordinates for any regions
118 where SNVs should be excluded from analysis.

119 **Architecture**

120 Execution of SNVPhyl begins with the “Repeat Identification” stage. This stage identifies internal repeat
121 regions on the reference genome using MUMMer (v3.23) (24) and generates a masking file containing the
122 locations of repetitive regions to exclude from analysis. This file is concatenated to the user-supplied
123 masking file, if defined, and used in later analysis stages.

124 The “Mapping/Variant Calling” stage (detailed in Figure 1.b) aligns the supplied reads to the
125 reference genome using the appropriate mapping mode (paired-end or single-end). Reference mapping is
126 performed using SMALT (v.0.7.5) (<http://www.sanger.ac.uk/science/tools/smalt-0>), which outputs a read
127 pileup. In the “Mapping Quality” stage, SNVPhyl evaluates each pileup for the mean coverage across a
128 user-defined proportion of the reference genome (e.g., 10X coverage across at least 80% of the genome).
129 Any sequenced genomes that do not meet the minimum mean coverage threshold are flagged for further
130 assessment.

131 The variant calling stages of SNVPhyl use two independent variant callers, FreeBayes (version
132 0.9.20) (25), and the SAMtools and BCFtools packages (26, 27). FreeBayes is run using the haploid
133 variant calling mode and the resulting variants are filtered to remove insertions/deletions and split
134 complex variant calls. SAMtools and BCFtools are run independently of FreeBayes and are used to
135 confirm the FreeBayes variant calls and generate base calls for non-variant positions.

136 The “Variant Consolidation” stage combines both sets of variant and non-variant calls into a
137 merged file, flagging mismatches between variant callers. Base calls below the defined minimum read
138 coverage are identified and flagged. The merged base calls are scanned for positions that do not pass the
139 minimum SNV abundance ratio (ratio of reads supporting the SNV with respect to the depth of coverage
140 at a site) and minimum mean mapping quality. These base calls are removed from the merged base calls
141 file. The remaining base calls that pass all these criteria are defined as either a high quality SNV (hqSNV)

142 or a high quality non-variant base call. The hqSNVs are optionally scanned to identify high density SNV
143 regions. These regions are identified by passing a sliding window of a given size along the genome and
144 counting the number of SNVs within the window that exceed a given SNV density threshold. The high-
145 density SNV regions are recorded in a tab-delimited file and used to mask potential recombinant regions.

146 The “SNV Alignment Generation” stage examines the merged base calls to generate a table of
147 identified variants and an alignment of hqSNVs and high-quality non-variant bases. The hqSNVs are
148 evaluated and assigned a status using the base calls at the same reference genome position for every
149 isolate. A status of “valid” is assigned when the base calls from all isolates in the same position pass the
150 minimum criteria (hqSNVs or high-quality non-variants). These base calls are incorporated into the SNV
151 alignment used for phylogeny generation. A status of “filtered-coverage” is assigned when one or more
152 isolates fail the minimum coverage threshold and the failed isolates’ base calls are annotated as ‘-’
153 (indicating no nucleotide or a gap). A status of “filtered-mpileup” is assigned when one or more isolates
154 have conflicting base calls between FreeBayes and SAMtools/BCFtools and the conflicting isolates’ base
155 calls are annotated as ‘N’ (indicating any nucleotide nonspecifically). A status of “filtered-invalid” is
156 assigned when the identified hqSNV overlaps one of the masked locations. The hqSNVs, base calls, and
157 assigned status are recorded in the SNV table and saved for later inspection. The SNV table can be used
158 to re-generate the downstream SNV alignment and phylogenetic tree without re-running the
159 computationally intensive reference mapping and variant calling steps.

160 **Output**

161 The final phylogeny is generated using the SNV alignment consisting of hqSNVs with a “valid” status.
162 This alignment is run through PhyML (28) with the GTR+ γ model as default and tree support values
163 estimated using PhyML’s approximate likelihood ratio test (29). The SNV alignment is also used to
164 generate an all-against-all SNV distance matrix. This matrix lists the pair-wise distances between every
165 isolate, using only the “valid” hqSNVs.

166 Additional files are provided to assist in evaluating the quality of the SNVPhyl analysis. The
167 “SNV Filter Stats” stage summarizes the quality and counts of the identified SNVs. The “SNV Alignment
168 Generation” stage summarizes the proportion of the reference genome passing all the necessary filters for
169 every isolate—the (non-masked) core genome—as well as the portion of the genome failing any filter or
170 excluded by the masking file.

171 **Simulated data**

172 We evaluated SNVPhyl’s sensitivity and specificity for SNV identification using simulated mutations
173 derived from a reference genome. The closed and finished *E. coli* str. Sakai (NC_002695) along with the
174 two plasmids (NC_002128 and NC_002127) was chosen as the reference genome (combined length of
175 5,594,477 bp). We constructed a variant genome by randomly mutating 10,000 base locations on the
176 reference genome. We repeated the procedure, using the same 10,000 base locations but different
177 mutations, to generate a total of three variant genomes. We included the unmodified reference genome in
178 the test set to serve as a positive control. The simulated variants for each genome were recorded in a table
179 for later comparisons. The constructed genomes were run through art_illumina (version
180 ChocolateCherryCake) (30) to generate paired-end reads with 2x250 bp length and 30X mean coverage.
181 The resultant reads along with the reference genome were run through SNVPhyl with repeat masking
182 enabled but with no SNV density filtering.

183 The SNV table produced by SNVPhyl was compared to the table of simulated variants to
184 determine their sensitivity and specificity. We define a true positive (TP) as a matching row in both
185 variant tables where both the position as well as base calls for each simulated genome is identical. A
186 variant detected by SNVPhyl not matching the criteria for a true positive is a false positive (FP). A true
187 negative (TN) is defined as all non-variant positions that were excluded by SNVPhyl. A false negative
188 (FN) is defined as a row in the simulated variant table where either the position or a base call did not
189 match any corresponding entry in the table of detected variants by SNVPhyl. Using these definitions,
190 sensitivity is calculated as $TP/(TP + FN)$ while specificity is calculated as $TN/(TN + FP)$.

191 **SNV density filtering evaluation**

192 We evaluated SNVPhyl's ability to mask recombination by comparing the resultant phylogenetic trees
193 and identified SNVs to those detected and removed by the recombination detection software package
194 Gubbins (18). Our test data consisted of 11 *Streptococcus pneumoniae* genomes along with the reference
195 genome ATCC 700669 (FM211187) that had previously been published (31) and made available as
196 sequence reads on NCBI (Table S1) and as a whole-genome alignment (the PMEN1 dataset from
197 <https://sanger-pathogens.github.io/gubbins/>). We downloaded this alignment, appended the reference
198 genome, and processed the resulting file through Gubbins to identify and mask recombinant SNVs. The
199 identified SNVs were filtered to remove gaps and masked recombination ('-' and 'N' characters) and the
200 resulting SNVs we defined as the "truth" set used to generate the true/false positive/negative values—
201 defined as for the "Simulated data" section. These Gubbins-identified SNVs were also used to construct a
202 phylogenetic tree with PhyML and compared with SNVPhyl's phylogenetic trees numerically using K
203 tree scores (32) and visually using phytools (33). K tree scores allow for similarity comparisons of many
204 phylogenetic trees against a single reference tree. Each tree is re-scaled by a factor, K, based on the
205 reference tree size and a score is produced taking into account differences in both topology and branch
206 lengths. Comparing the scores of all trees provides a measure of similarity to the reference tree, with more
207 similar trees producing a score closer to 0.

208 We downloaded sequence reads for the test dataset from NCBI, identifying and combining
209 multiple sequencing runs for each strain to a single set of sequence reads with the help of SRADB (34).
210 Using the combined sequence reads we ran SNVPhyl under a number of scenarios. For each scenario we
211 compared the SNVs and phylogenetic trees to the "truth" dataset described above. In the first run, we
212 performed no SNV density filtering. For all subsequent runs we adjusted the density-filtering parameters
213 to remove SNVs occurring at a density of 2 or more within a moving window of 20, 100, 500, 1000, and
214 2000 bp. We evaluated an additional scenario using a combination of SNVPhyl and Gubbins for
215 recombination masking. We ran SNVPhyl with no SNV density filtering and incorporated the identified

216 variants into the reference genome to generate a whole genome alignment. The whole-genome alignment
217 was processed with Gubbins to identify non-recombinant SNVs and to construct a phylogenetic tree.

218 **Parameter optimization**

219 We evaluated SNVPhyl's parameter settings and resulting accuracy at differentiating outbreak isolates
220 using a set of 59 sequenced and published *Salmonella enterica* serovar Heidelberg genomes (35), which
221 were previously deposited in the NCBI Sequence Read Archive (Table S2). We chose this dataset as it
222 contained sequence data for strains from several unrelated outbreaks—referred to as “outbreak 1”,
223 “outbreak 2”, and “outbreak 3”—along with additional background strains, allowing us to evaluate
224 SNVPhyl's ability to differentiate the outbreak strains under different scenarios. Sequence read data was
225 subsampled with seqtk (<https://github.com/lh3/seqtk>) such that the genome with the least amount of
226 sequence data, SH12-006, was set to 30X mean coverage (calculated as: mean coverage = count of base
227 pairs in all reads / length of reference genome). Other genomes were subsampled to maintain their relative
228 proportion of mean read coverage to SH12-006. *Salmonella* Heidelberg str. SL476 (NC_011083) was
229 selected as the reference genome. We optimized the SNVPhyl parameters for this dataset according to the
230 following four scenarios: 1) adjusting the minimum base coverage parameter used to call a variant while
231 keeping the number of reads in the dataset fixed; 2) subsampling the reads of a single WGS sample at
232 different mean coverage levels while keeping the minimum base coverage parameter fixed; 3) adjusting
233 the minimum SNV abundance ratio for calling a variant; and 4) adjusting the amount of contamination in
234 the dataset to determine its effect on variant calling accuracy.

235 In the first scenario we ran the SNVPhyl pipeline using the default parameters except for the
236 minimum base coverage, which was adjusted to 5X, 10X, 15X, and 20X. In the second scenario we kept
237 the minimum base coverage parameter fixed at 10X, while one of the samples (SH13-001) was
238 subsampled to mean sequencing coverages of 30X, 20X, 15X, and 10X. In the third scenario the
239 minimum SNV abundance ratio was adjusted to 0.25, 0.5, 0.75, and 0.9. In the fourth scenario, a sample
240 from “outbreak 2” (SH13-001 with mean coverage 71X) was chosen as a candidate for simulating

241 contamination. A sample from the unrelated “outbreak 1” (SH12-001) was selected as the source of
242 contaminant reads. The reads were subsampled and combined such that SH13-001 (“outbreak 2”)
243 remained at 71X mean coverage but was contaminated with reads from SH12-001 (“outbreak 1”) at 5%,
244 10%, 20%, and 30%. All samples were run through SNVPhyl for each of these contamination ratios.

245 The phylogenetic trees produced by SNVPhyl were evaluated for concordance with the outbreak
246 epidemiological data using the following criteria: 1) all outbreak isolates group monophyletically, and 2)
247 the SNV distance between any two isolates within an outbreak clade is less than 5 SNVs, a number
248 identified in the previous study (35) as the maximum SNV distance between epidemiologically related
249 samples within these particular outbreaks. Both conditions were tested using the APE package within R
250 (36).

251 **Results**

252 **Validation against simulated data**

253 We measured SNVPhyl’s sensitivity and specificity by introducing random mutations along the *E. coli*
254 Sakai reference genome and compared these mutations with those detected by SNVPhyl (Table 2). Of the
255 10,000 mutated positions introduced, SNVPhyl reported 9,116 true positives and 0 false positives
256 resulting in a sensitivity and specificity of 0.91 and 1.0.

257 Positions on the reference genome that contain a low-quality base call, or exist in repetitive
258 regions are excluded from downstream analysis by SNVPhyl. However, lower-quality variant-containing
259 sites along with variants in repetitive regions are saved by SNVPhyl in the variant table with a “filtered”
260 status. Comparing the combination of low-quality and high-quality variant sites detected by SNVPhyl, we
261 found 9,573 true positives and 51 false positives resulting in a sensitivity and specificity of 0.96 and 1.0
262 (after rounding). Of the 51 false positives, 48 were considered as false positives due to insufficient read
263 coverage in one of the samples to call a high quality variant, thus resulting in a call of a gap (‘-’) as

264 opposed to the true base call. Only 3 of the false positives were a result of miscalled bases with sufficient
265 read coverage, and these occurred in repetitive regions of the genome.

266 **SNV density filtering evaluation**

267 We compared SNVPhyl's density filtering against the Gubbins software for detection and removal of
268 recombination in a collection of WGS reads from 11 *Streptococcus pneumoniae* genomes along with the
269 reference genome ATCC 700669 (Table 3, Figure S1). With no SNV density filtering SNVPhyl properly
270 identified 142 SNV-containing sites (true positives) but included 2,159 additional SNV sites (false
271 positives). These false positives skew the resulting phylogenetic tree by increasing the length of one of
272 the branches. The phylogenetic tree is compared with the tree produced with Gubbins, resulting in a K
273 tree score of 0.419.

274 We reanalyzed the dataset with high-density SNV masking enabled, using a range of variant
275 density cutoffs. We found the density-filtering criteria of 2 SNVs in a 500 bp window and 2 SNVs in a
276 1000 bp window performed near-equally in producing a phylogenetic tree resembling the tree produce by
277 Gubbins based on the K tree scores of 0.045 and 0.044, both much lower than the score of 0.419 for no
278 SNV density filtering. With these filtering criteria, SNVPhyl identified 133 true positives and 12 false
279 positives (for 2 SNVs in 500 bp) and 125 true positives and 6 false positives (for 2 SNVs in 1000 bp).

280 We also investigated the effect of generating a whole-genome alignment—by incorporating
281 SNVPhyl-identified variants without SNV density filtering into the reference genome—for a more
282 thorough analysis with the recombination-detection software Gubbins. We were able to identify 138 true
283 positives in the alignment at the expense of 10 false positives and a K tree score of 0.037, a result closely
284 matching the use SNVPhyl's density filtering criteria.

285 **Parameter optimization**

286 We evaluated SNVPhyl's capability to differentiate between epidemiologically related and unrelated
287 samples using a WGS dataset consisting of 59 *Salmonella enterica* serovar Heidelberg genomes from

288 three unrelated outbreaks. We ran SNVPhyl with this data under a number of scenarios: 1) varying the
289 minimum base coverage required by SNVPhyl to call a variant, 2) subsampling the reads of an individual
290 bacterial sample, 3) varying the minimum SNV abundance ratio, and 4) testing the ability to generate
291 accurate phylogenetic trees in the presence of contamination. We tested the SNVPhyl results for
292 phylogenetic concordance to epidemiological data (Table 4, Figure S2).

293 For the first scenario we found that as the minimum coverage threshold was increased, the
294 percent of the reference genome identified as part of the core genome and number of SNV-containing
295 sites was reduced (from 95% core and 317 SNVs to 54% core and 165 SNVs). At 15X minimum
296 coverage (81% core and 262 SNVs) and lower all three outbreaks grouped into monophyletic clades.
297 Failure occurred at a minimum coverage of 20X (54% core and 165 SNVs), where the outbreak 2 isolates
298 failed to constitute a separate clade.

299 For the second scenario, one of the samples was subsampled to reduce the mean coverage relative
300 to all other samples while keeping the minimum coverage parameter of 10X in SNVPhyl fixed. At a
301 mean coverage of 15X (with 242 SNVs identified and 76% core) SNVPhyl grouped all three outbreaks
302 into monophyletic clades. However, at a lower mean coverage of 10X (155 SNVs and 47% core)
303 SNVPhyl failed to group one of the outbreaks into a monophyletic clade. Similar to the first scenario, the
304 percentage of the reference genome considered core as well as the number of SNVs identified was
305 reduced as the mean coverage of one of the samples was lowered.

306 For the third scenario, the SNV abundance ratio—defining the ratio of SNV-supporting bases
307 needed to identify a variant as high-quality—was adjusted incrementally. Each set of outbreak isolates
308 grouped into a clade with a maximum SNV distance less than 5 SNVs above a ratio of 0.5. At a ratio of
309 0.5 the maximum SNV distance within outbreak 2 was exactly 5 SNVs while for a ratio of 0.25 the
310 maximum SNV distance in outbreak 2 was 44 SNVs. The percentage of the reference genome identified
311 as part of the core genome remained the same at 92%.

312 For the fourth scenario, we examined the robustness of SNVPhyl to cross-contamination of
313 closely related samples. Current methods of contamination detection often focus on taxonomic
314 classification of genomic content (37). However, contamination by closely related isolates can go
315 undetected. We simulated contamination for an isolate in outbreak 2 by an isolate in outbreak 1. We
316 found that SNVPhyl was able to accurately differentiate all three outbreaks with up to 10% read
317 contamination; however the number of SNVs dropped from 298 SNVs at 5% contamination, to 260 SNVs
318 at 20% contamination, where the failure was due to removal of the majority of unique SNVs that
319 differentiated outbreak 1 from the background isolates.

320 Discussion

321 The availability of WGS data from microbial genomes represents a tremendous opportunity for infectious
322 disease surveillance and outbreak response. Emerging analytical methods, such as gene-by-gene or SNV-
323 based, require that bioinformatics pipelines be designed with usability by non-bioinformaticians in mind
324 and which can be easily incorporated into existing systems. An overview of current phylogenomic
325 methods appears in (38) and a comparison of SNVPhyl's design with that of other popular pipelines
326 appears in Table 1 (a detailed investigation comparing the performance of SNVPhyl with other pipelines
327 is the subject of a forthcoming manuscript). We designed SNVPhyl to be both flexible and scalable in its
328 usage in order to meet the needs and abilities of most labs. SNVPhyl gains much of this flexibility
329 through its implementation as a Galaxy workflow, which enables execution in environments from single
330 machines to high-scale computer clusters, from third-party web-based environments to local installations.
331 Galaxy provides a user-friendly interface but also provides an API, which is used to implement a
332 command-line interface for SNVPhyl. The SNVPhyl pipeline is also integrated within the IRIDA
333 platform (<http://irida.ca>), which provides an integrated “push-button” system for genomic epidemiology.
334 However, implementing SNVPhyl through Galaxy has some disadvantages. Notably, Galaxy is more
335 complex and so more cumbersome to install than a simpler command-line based pipeline. To address this

336 we have made SNVPhyl available as simple to install virtual machine and Docker images, although these
337 options cannot be straightforwardly implemented in a high performance computing environment.

338 Several factors can influence the ability to accurately call SNVs when using a reference mapping
339 approach (39). As well, there are aspects of the datasets—such as recombination and population
340 diversity—that can influence the phylogenetic analysis of identified SNVs. To assist in selecting proper
341 parameters for SNVPhyl and gauging performance on different datasets we have assessed SNVPhyl under
342 a variety of situations: SNV calling accuracy with simulated data, recombination masking, and the ability
343 to differentiate outbreak isolates from non-outbreak isolates under differing parameters and data qualities.

344 Our assessment of SNV calling accuracy shows that SNVPhyl can detect SNVs and produce a
345 SNV alignment with high sensitivity and specificity (Table 2). Of the variants that went undetected by
346 SNVPhyl a large proportion were due to the quality thresholds and masking procedures implemented by
347 SNVPhyl to remove incorrectly called or problematic SNVs (e.g., SNVs in internal repeats on the
348 reference genome). While these quality procedures generate many false negatives they also eliminate
349 many false positive variants—a reduction of 51 to 0 false positives at a cost of an additional 457 false
350 negatives in the simulated dataset. However, all detected variation across all genomes is recorded in a
351 table produced by SNVPhyl and additional software is provided for more detailed analysis of these
352 variants.

353 Phylogenetics assumes descent with modification, but recombination (horizontal gene transfer)
354 violates this assumption and its presence can confound the resulting phylogeny leading to
355 misinterpretations on the clonal relationship of isolates (40). Recombination detection software exists and
356 can be used for the construction of phylogenetic trees based on vertically inherited information (18, 19,
357 41). However, these programs require the pre-construction of whole genome alignments and can only be
358 run on a single machine, which limits their utility for routine application to large collections of WGS
359 reads.

360 SNVPhyl implements a basic but rapid method for detection and masking of recombinant sites by
361 searching for SNV-dense regions above a defined density in a sliding window. We evaluated SNVPhyl's
362 recombination-masking method in comparison to the Gubbins software package which was run on a
363 previously generated whole-genome alignment (Table 3, Figure S1). We found that SNVPhyl removes
364 the majority of recombinant SNVs (from 2,159 SNVs with no recombination masking to 6 SNVs when
365 masking regions with 2 SNVs in a 1000 bp window). However, SNVPhyl also removes some non-
366 recombinant SNVs (reduced from 142 SNVs with no masking to 125 SNVs with 2 SNVs in 1000 bp).
367 Removal of a greater number of recombinant SNVs is possible by increasing the window size, but this
368 removes additional non-recombinant SNVs and reduces the information available in the phylogenetic tree
369 and so concordance with other recombination-masking procedures (based on K tree scores).

370 SNVPhyl's method of detecting high-density SNV regions can be executed independently for
371 each genome. Independent execution is easily distributed across multiple nodes within a compute cluster,
372 enhancing the scalability over large datasets. However, SNVPhyl requires the SNV density to be set *a*
373 *priori* and may not be appropriate for organisms with complex evolutionary dynamics or for genome
374 sequences from organisms spanning a large phylogenetic distance. We suspect that the optimal
375 parameters will vary based on the particular organism under study and we would caution against relying
376 on default settings without further evaluation. SNVPhyl does not aim to be a rigorous recombination
377 detection and removal software package. However, SNVPhyl provides output files recording all the SNVs
378 detected, which can be used for further analysis if needed. In particular, additional tools are provided that
379 can produce a whole genome alignment correctly formatted for input into software such as Gubbins for a
380 thorough detection of recombination and construction of a phylogenetic tree from non-recombinant
381 SNVs.

382 A proper interpretation of the produced phylogenetic trees and SNV distances for associating
383 closely-related isolates requires knowledge of when to trust the results and when additional parameter or
384 data adjustments are necessary. To assist in defining these criteria we evaluated the performance of

385 SNVPhyl at clearly delineating different outbreak clades across four different scenarios (Table 4, Figure
386 S2).

387 In both the first and second scenarios we examined the effect of sequencing coverage on
388 identifying enough SNVs to properly differentiate outbreak isolates. In the first scenario, we adjusted the
389 minimum base coverage required to call a SNV from 5X to 20X without any additional subsampling of
390 reads. We found that SNVPhyl succeeded in differentiating outbreak isolates at coverages up to 15X, but
391 at a minimum coverage of 20X SNVPhyl failed to differentiate the outbreak isolates due to removal of
392 too many SNVs (from 317 SNVs to 165 SNVs). In the second scenario, we subsampled one of the
393 isolates along the mean read coverage values from 30X to 10X while keeping the minimum base coverage
394 parameter in SNVPhyl fixed at 10X. We found SNVPhyl succeeded in differentiating outbreak isolates at
395 a mean coverage of 15X and above, but failed to differentiate outbreak isolates at a mean coverage of
396 10X due to removal of too many SNVs (reduced from 299 SNVs to 155 SNVs). Both cases show that a
397 high base coverage threshold for variant calling relative to the mean coverage of the lowest sample leads
398 to falsely identifying samples as being related due to removal of too many SNVs (20X minimum
399 coverage / 30X lowest sample mean coverage for failure in the first scenario, and 10X minimum coverage
400 / 10X lowest sample mean coverage for failure in the second scenario). However, a high minimum base
401 coverage threshold or too little sequencing data can be detected by examining the percentage of the
402 reference genome considered as part of the core genome by SNVPhyl. A low value can indicate either a
403 poorly related reference genome, or that large portions of the genomes are removed from the analysis
404 (drop from 95% to 54% in the first scenario and 92% to 47% in the second scenario). We would
405 recommend searching for such low values in the percent core to gauge whether or not base coverage (or
406 possibly reference genome selection) is an issue for the SNVPhyl results.

407 In the third scenario (Table 4, Figure S2.c) we adjusted the SNV abundance ratio among values
408 from 0.25 to 0.9. We found that SNVPhyl successfully differentiated outbreak isolates above a ratio of
409 0.5, but at a ratio of 0.5 the maximum SNV distance between isolates within an outbreak exceeded our

410 threshold of less than 5 SNVs. However, unlike the minimum base coverage value, the percent of the
411 reference genome identified as the core genome remained the same (92%). We recommend keeping this
412 setting fixed at a higher value, with the default set at 0.75.

413 In the fourth scenario we simulated contamination between two closely-related isolates from two
414 different outbreaks by mixing reads at differing proportions (Table 4, Figure S2.d). Our findings indicate
415 that SNVPhyl is able to handle low amounts of mixed sample contamination (up to 10%). A higher
416 proportion of contaminated reads can lead to removal of SNVs due to not meeting quality thresholds
417 (from 298 SNVs with 5% contamination to 260 SNVs at 20% contamination where failure occurred) and
418 so incorrectly implying relatedness between samples. Similar to the third scenario, the percentage of the
419 reference genome identified as the core genome remained fixed at 92%. While SNVPhyl is able to
420 differentiate outbreak isolates at low levels of contamination SNVPhyl cannot be used to evaluate the
421 degree of contamination. Thus, we would not recommend the straightforward application of SNVPhyl to
422 contaminated datasets without further assessment of the degree of contamination, either through
423 taxonomic identification software such as Kraken (42); or, for closely-related isolates, through inspection
424 of the variant calling and read pileup information provided by SNVPhyl.

425 Our analysis suggests that great care must be taken to reduce sources of noise in genome-wide
426 SNV analysis. Some of this noise relates to quality thresholds for calling high quality SNVs, of which a
427 careful balance is required to eliminate false positives without removal of too many true variants. Other
428 sources include aspects of the WGS datasets or organisms under study such as the presence of
429 contamination or recombination. The studied cases highlight how SNVPhyl is able to produce accurate
430 phylogenetic trees under a wide variety of data qualities, and demonstrate how to detect inaccurate trees
431 using additional information generated by SNVPhyl.

432 **Conclusion**

433 SNVPhyl provides an easy-to-use pipeline for processing whole genome sequence reads to identify SNVs
434 and produce a phylogenetic tree. We have shown that SNVPhyl is capable of producing accurate results
435 on even very closely related bacterial isolates under a wide variety of parameter settings and sequencing
436 data qualities. SNVPhyl is distributed as a pipeline within Galaxy and is integrated within the IRIDA
437 platform, providing a “push button” system for generating whole genome phylogenies within a larger
438 WGS data management and genomic epidemiology system designed for use in clinical, public health, and
439 food regulatory environments.

440 **Author Contributions**

441 A.P., N.K., C.N., A.R., E.T., R.B., W.H., F.B., M.G., and G.V.D. wrote the manuscript. A.P., L.K., A.R.,
442 and G.V.D. contributed to the initial software design while A.P., P.M., C.S., N.K., A.R., K.W., R.Z., and
443 G.V.D. extended the design to a fully automated pipeline. A.P., P.M., C.S., and G.V.D. wrote the
444 software with L.K., J.C., Mari. I., and Mark I. providing additional scripts or wrapping Galaxy tools.
445 A.P., C.S., A.R., K.W., and G.V.D. designed the validation experiments. A.P. and C.S. performed the
446 validation experiments. A.P. and G.V.D. interpreted the results.

447 **Acknowledgments**

448 The authors would like to acknowledge Cheryl Tarr for initial inspiration and contribution to the design of
449 SNVPhyl as well as Lauren Sluskey and Brian Yeo for their contributions during development of the
450 pipeline. The authors would also like to acknowledge the Galaxy team and Galaxy community for their
451 rapid response to issues and feature requests during the development of SNVPhyl as well as the
452 integration of some bioinformatics tools within Galaxy that are used by SNVPhyl.

453 **Funding**

454 This work was supported by the Genomics Research and Development Initiative, Genome Canada, and
455 Genome British Columbia.

456 **References**

457

- 458 1. Hendriksen RS, Price LB, Schupp JM, Gillece JD, Kaas RS, Engelthaler DM, et al. Population genetics
459 of *Vibrio cholerae* from Nepal in 2010: evidence on the origin of the Haitian outbreak. *MBio*. 2011 Sep
460 1;2(4):e00157-11.
- 461 2. Katz LS, Petkau A, Beaulaurier J, Tyler S, Antonova ES, Turnsek MA, et al. Evolutionary dynamics of
462 *Vibrio cholerae* O1 following a single-source introduction to Haiti. *MBio*. 2013 Jul
463 2;4(4):10.1128/mBio.00398-13.
- 464 3. Frerichs RR, Keim PS, Barraix R, Piarroux R. Nepalese origin of cholera epidemic in Haiti. *Clin*
465 *Microbiol Infect*. 2012 Jun;18(6):E158-63.
- 466 4. Gardy JL, Johnston JC, Sui SJH, Cook VJ, Shah L, Brodtkin E, et al. Whole-genome sequencing and
467 social-network analysis of a tuberculosis outbreak. *N Engl J Med*. 2011;364(8):730-9.
- 468 5. Roetzer A, Diel R, Kohl TA, Rückert C, Nübel U, Blom J, et al. Whole genome sequencing versus
469 traditional genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a longitudinal
470 molecular epidemiological study. *PLoS Med*. 2013;10(2):e1001387.
- 471 6. Holmes A, Allison L, Ward M, Dallman TJ, Clark R, Fawkes A, et al. Utility of Whole-Genome
472 Sequencing of *Escherichia coli* O157 for Outbreak Detection and Epidemiological Surveillance. *J Clin*
473 *Microbiol*. 2015 Nov;53(11):3565-73.
- 474 7. Sánchez-Busó L, Comas I, Jorques G, González-Candelas F. Recombination drives genome evolution
475 in outbreak-related *Legionella pneumophila* isolates. *Nat Genet*. 2014;46(11):1205-11.
- 476 8. Allard MW, Strain E, Melka D, Bunning K, Musser SM, Brown EW, et al. Practical Value of Food
477 Pathogen Traceability through Building a Whole-Genome Sequencing Network and Database. *J Clin*
478 *Microbiol*. 2016 Aug;54(8):1975-83.
- 479 9. Franz E, Gras LM, Dallman T. Significance of whole genome sequencing for surveillance, source
480 attribution and microbial risk assessment of foodborne pathogens. *Current Opinion in Food Science*.
481 2016;8:74-9.
- 482 10. Ashton PM, Nair S, Peters TM, Bale JA, Powell DG, Painset A, et al. Identification of *Salmonella* for
483 public health surveillance using whole genome sequencing. *PeerJ*. 2016 Apr 5;4:e1752.

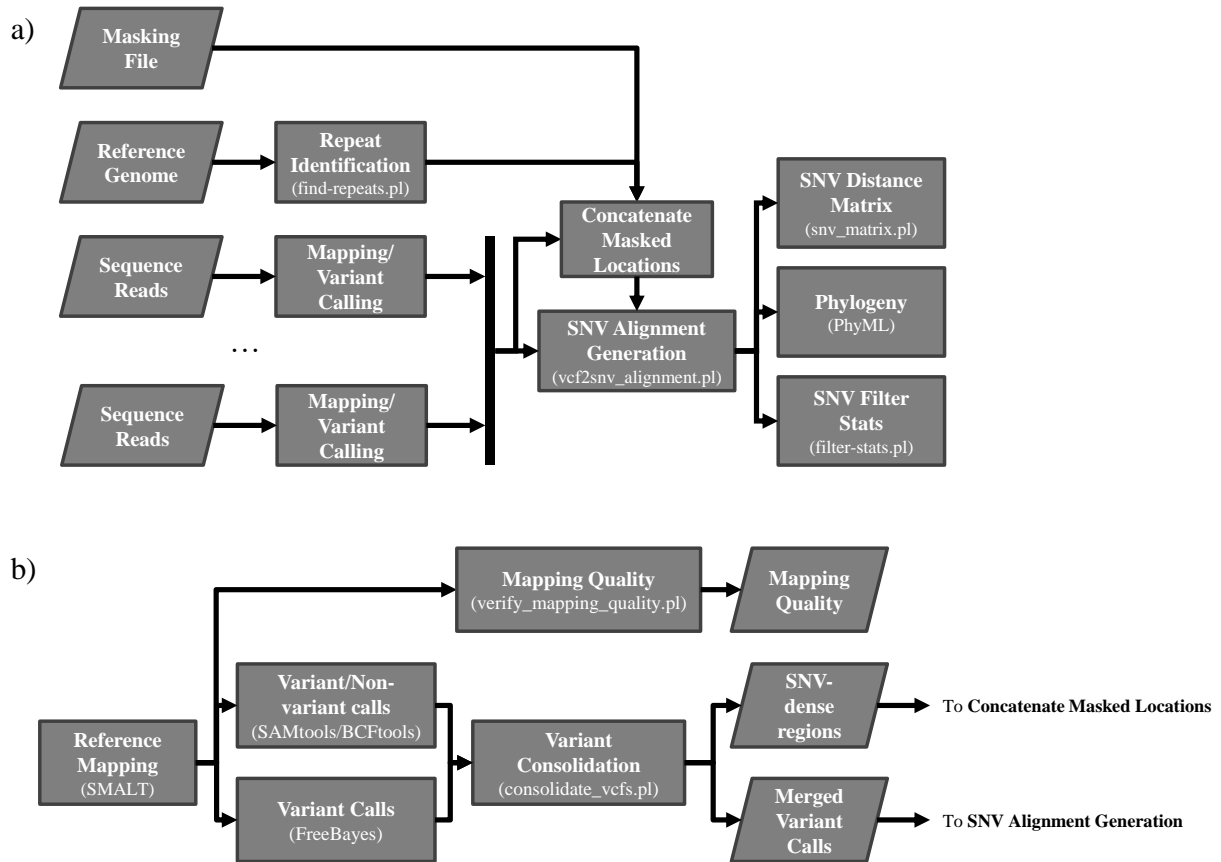
- 484 11. Maiden MC, Jansen van Rensburg MJ, Bray JE, Earle SG, Ford SA, Jolley KA, et al. MLST revisited:
485 the gene-by-gene approach to bacterial genomics. *Nat Rev Microbiol*. 2013 Oct;11(10):728-36.
- 486 12. Moura A, Criscuolo A, Pouseele H, Maury MM, Leclercq A, Tarr C, et al. Whole genome-based
487 population biology and epidemiological surveillance of *Listeria monocytogenes*. *Nature Microbiology*.
488 2016 10/10;2:16185.
- 489 13. Kwong JC, Mercoulia K, Tomita T, Easton M, Li HY, Bulach DM, et al. Prospective Whole-Genome
490 Sequencing Enhances National Surveillance of *Listeria monocytogenes*. *J Clin Microbiol*. 2016
491 Feb;54(2):333-42.
- 492 14. Bertels F, Silander OK, Pachkov M, Rainey PB, van Nimwegen E. Automated reconstruction of
493 whole-genome phylogenies from short-sequence reads. *Mol Biol Evol*. 2014 May;31(5):1077-88.
- 494 15. Jackson BR, Tarr C, Strain E, Jackson KA, Conrad A, Carleton H, et al. Implementation of
495 Nationwide Real-time Whole-genome Sequencing to Enhance Listeriosis Outbreak Detection and
496 Investigation. *Clinical Infectious Diseases*. 2016 April 18.
- 497 16. Kaas RS, Leekitcharoenphon P, Aarestrup FM, Lund O. Solving the problem of comparing whole
498 bacterial genomes across different sequencing platforms. *PLoS One*. 2014 Aug 11;9(8):e104984.
- 499 17. Davis S, Pettengill JB, Luo Y, Payne J, Shpuntoff A, Rand H, et al. CFSAN SNP Pipeline: an
500 automated method for constructing SNP matrices from next-generation sequence data. *PeerJ Computer
501 Science*. 2015;1:e20.
- 502 18. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid phylogenetic
503 analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic
504 Acids Res*. 2015 Feb 18;43(3):e15.
- 505 19. Didelot X, Wilson DJ. ClonalFrameML: efficient inference of recombination in whole bacterial
506 genomes. *PLoS Comput Biol*. 2015 Feb 12;11(2):e1004041.
- 507 20. Sahl JW, Lemmer D, Travis J, Schupp JM, Gillece JD, Aziz M, et al. NASP: an accurate, rapid
508 method for the identification of SNPs in WGS datasets that supports flexible input and output formats.
509 *Microbial Genomics*. 2016;2(8).
- 510 21. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Cech M, et al. The Galaxy platform
511 for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res*. 2016
512 Jul 8;44(W1):W3-W10.
- 513 22. Afgan E, Sloggett C, Goonasekera N, Makunin I, Benson D, Crowe M, et al. Genomics Virtual
514 Laboratory: A Practical Bioinformatics Workbench for the Cloud. *PLoS One*. 2015 Oct
515 26;10(10):e0140829.
- 516 23. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, et al. Dissemination of
517 scientific software with Galaxy ToolShed. *Genome Biol*. 2014 Feb 20;15(2):403.
- 518 24. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open
519 software for comparing large genomes. *Genome Biol*. 2004;5(2):R12.

- 520 25. E. G, G. M. Haplotype-based variant detection from short-read sequencing. ArXiv e-prints. 2012 jul.
- 521 26. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
522 format and SAMtools. *Bioinformatics*. 2009 Aug 15;25(16):2078-9.
- 523 27. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and
524 population genetical parameter estimation from sequencing data. *Bioinformatics*. 2011 Nov
525 1;27(21):2987-93.
- 526 28. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by
527 maximum likelihood. *Syst Biol*. 2003 Oct;52(5):696-704.
- 528 29. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and
529 methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst*
530 *Biol*. 2010 May;59(3):307-21.
- 531 30. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator.
532 *Bioinformatics*. 2012 Feb 15;28(4):593-4.
- 533 31. Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J, van der Linden M, et al. Rapid pneumococcal
534 evolution in response to clinical interventions. *Science*. 2011 Jan 28;331(6016):430-4.
- 535 32. Soria-Carrasco V, Talavera G, Igea J, Castresana J. The K tree score: quantification of differences in
536 the relative branch length and topology of phylogenetic trees. *Bioinformatics*. 2007 Nov 1;23(21):2954-6.
- 537 33. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other things). *Methods*
538 *in Ecology and Evolution*. 2012;3:217-23.
- 539 34. Zhu Y, Stephens RM, Meltzer PS, Davis SR. SRADB: query and use public next-generation
540 sequencing data from within R. *BMC Bioinformatics*. 2013 Jan 17;14:19,2105-14-19.
- 541 35. Bekal S, Berry C, Reimer AR, Van Domselaar G, Beaudry G, Fournier E, et al. Usefulness of High-
542 Quality Core Genome Single-Nucleotide Variant Analysis for Subtyping the Highly Clonal and the Most
543 Prevalent *Salmonella enterica* Serovar Heidelberg Clone in the Context of Outbreak Investigations. *J Clin*
544 *Microbiol*. 2016 Feb;54(2):289-95.
- 545 36. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language.
546 *Bioinformatics*. 2004 Jan 22;20(2):289-90.
- 547 37. Koren S, Treangen TJ, Hill CM, Pop M, Phillippy AM. Automated ensemble assembly and validation
548 of microbial genomes. *BMC Bioinformatics*. 2014 May 3;15:126,2105-15-126.
- 549 38. Lynch T, Petkau A, Knox N, Graham M, Van Domselaar G. A Primer on Infectious Disease Bacterial
550 Genomics. *Clinical Microbiology Reviews*. 2016 October 01;29(4):881-913.
- 551 39. Olson ND, Lund SP, Colman RE, Foster JT, Sahl JW, Schupp JM, et al. Best practices for evaluating
552 single nucleotide variant calling methods for microbial genomics. *Front Genet*. 2015 Jul 7;6:235.

- 553 40. Croucher NJ, Harris SR, Grad YH, Hanage WP. Bacterial genomes in epidemiology--present and
554 future. *Philos Trans R Soc Lond B Biol Sci*. 2013 Feb 4;368(1614):20120202.
- 555 41. Marttinen P, Hanage WP, Croucher NJ, Connor TR, Harris SR, Bentley SD, et al. Detection of
556 recombination events in bacterial genomes from large population samples. *Nucleic Acids Res*. 2012
557 Jan;40(1):e6.
- 558 42. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
559 alignments. *Genome Biol*. 2014 Mar 3;15(3):R46,2014-15-3-r46.
- 560 43. Gardner SN, Slezak T, Hall BG. kSNP3.0: SNP detection and phylogenetic analysis of genomes
561 without genome alignment or reference genome. *Bioinformatics*. 2015 Sep 1;31(17):2877-8.
- 562 44. Treangen TJ, Ondov BD, Koren S, Phillippy AM. The Harvest suite for rapid core-genome alignment
563 and visualization of thousands of intraspecific microbial genomes. *Genome Biol*. 2014;15(11):524.
- 564 45. Ahmed SA, Lo C, Li P, Davenport KW, Chain PSG. From raw reads to trees: Whole genome SNP
565 phylogenetics across the tree of life. bioRxiv. 2015 Cold Spring Harbor Laboratory Press.
- 566

567 **Figures**

568 **Figure 1. a)** Overview of the SNVPhyl pipeline. Input to the pipeline is provided as a reference genome,
569 set of sequence reads for each isolate, and an optional list of positions to mask from the final results.
570 Repeat regions are identified on the reference genome and reference mapping followed by variant calling
571 is performed on the sequence reads. The resulting files are compiled together to construct a SNV
572 alignment and list of identified SNVs, which are further processed to construct a SNV distance matrix,
573 maximum likelihood phylogeny, and a summary of the identified SNVs. Individual software or scripts
574 are given in the parenthesis below each stage. **b)** An overview of the “Mapping/Variant Calling” stage of
575 SNVPhyl. Variants are called using two separate software packages and compiled together in the
576 “Variant Consolidation” stage. As output, a list of the validated variant calls, regions with high-density
577 SNVs, as well as quality information on the mean mapping coverage are produced and sent to further
578 stages.



579

580 Tables

581 **Table 1.** A comparison of whole-genome phylogenetic software.

Name	Input ¹	Parallel Computing ²	Distribution ³	Interface ⁴	Reference
CFSAN SNP Pipeline	sr	mn,mt	local	cl	(17)
CSIPhylogeny	sr,ag	n/a	web	gui	(16)
kSNP	sr,ag	mt	local	cl	(43)
Lyve-SET	sr,ag*	mn,mt	local	cl	https://github.com/lskatz/lyve-SET
NASP	sr,ag	mn,mt	local	cl	(20)
Parsnp	ag	mt	local	cl	(44)
PhaME	sr,ag	mt	local	cl	(45)
REALPHY	sr,ag	mt	web,local	gui,cl	(14)
Snippy	sr,ag*	mt	local	cl	https://github.com/tseemann/snippy
SNVPhyl	sr	mn,mt	local	gui,cl	http://snvphyl.readthedocs.io/

582 ¹sr = sequence reads, ag = assembled genome, ag* = assembled genome supported by generating

583 simulated reads

584 ²mn (multi-node) = provides capability to execute across multiple compute nodes, mt (multi-thread) =

585 provides multi-threading capability, n/a = not applicable (not locally installable)

586 ³local = locally distributed and installable software, web = software provided as a web service

587 ⁴cl = command-line interface, gui = graphical user interface

588 **Table 2.** SNV simulation results.

Comparison	Variant columns simulated	Non-variant columns	True positives	False positives	True negatives	False negatives	Specificity	Sensitivity
Valid SNVs ¹	10000	5584477	9116	0	5575361	884	1.0	0.91
All SNVs ²	10000	5584477	9573	51	5574853	427	1.0	0.96

589 ¹Valid SNVs: The number of SNV-containing sites detected that passed all thresholds to be considered

590 high quality for every isolate.

591 ²All SNVs: All the SNV-containing sites identified by SNVPhyl, including those where at least one

592 isolate did not have a high-quality base call or sites that were masked by the pipeline.

593 **Table 3.** A comparison of the SNVPhyl variant density filtering algorithm to the Gubbins system for

594 recombination detection.

Case	True Positives	False Positives	True Negatives	False Negatives	Sensitivity	Specificity	K tree score
No DF ¹	142	2159	2218849	23	0.861	0.999	0.419
2 in 20 ²	142	565	2220443	23	0.861	1.000	0.425
2 in 100 ²	142	155	2220853	23	0.861	1.000	0.377
2 in 500 ²	133	12	2221005	32	0.806	1.000	0.045
2 in 1000 ²	125	6	2221019	40	0.758	1.000	0.044
2 in 2000 ²	111	3	2221036	54	0.673	1.000	0.063
Gubbins/SNVPhyl ³	138	10	2221002	27	0.836	1.000	0.037

595 ¹No DF=A case of no SNV density filtering by SNVPhyl,

596 ²X in Y = Masking regions with a density of X variants in Y bases.

597 ³Gubbins/SNVPhyl = A whole genome alignment generated from SNVs identified by SNVPhyl and run
 598 through Gubbins.

599 **Table 4.** A comparison of the performance of SNVPhyl across a range of parameters and analysis
 600 scenarios.

#	Scenario	Parameters/Conditions	hqSNVs	% Core	Differentiated Outbreaks
1	Minimum Coverage	5X	317	95	Yes
		10X	301	92	Yes
		15X	262	81	Yes
		20X	165	54	No
2	Subsample Coverage Level	10X ¹	155	47	No
		15X ¹	242	76	Yes
		20X ¹	276	88	Yes
		30X ¹	299	92	Yes
3	SNV Abundance Ratio	0.25	351	92	No
		0.5	307	92	No
		0.75	301	92	Yes
		0.9	291	92	Yes
4	Contamination	5%	298	92	Yes
		10%	292	92	Yes
		20%	260	92	No
		30%	231	92	No

601 ¹These represent the mean coverage of one sample after subsampling reads and not the minimum
 602 coverage parameter of SNVPhyl (which is fixed at 10X).

603 **Supplementary Materials**

604 Supplementary materials are available in a separate file.