

SNVPhyl: a single nucleotide variant phylogenomics pipeline for microbial genomic epidemiology

Aaron Petkau,¹ Philip Mabon,¹ Cameron Sieffert,¹ Natalie C. Knox,¹ Jennifer Cabral,¹ Mariam Iskander,² Mark Iskander,² Kelly Weedmark,³ Rahat Zaheer,⁴ Lee S. Katz,⁵ Celine Nadon,¹ Aleisha Reimer,¹ Eduardo Taboada,¹ Robert G. Beiko,⁶ William Hsiao,⁷ Fiona Brinkman,⁸ Morag Graham¹ and Gary Van Domselaar^{1,*}

Abstract

The recent widespread application of whole-genome sequencing (WGS) for microbial disease investigations has spurred the development of new bioinformatics tools, including a notable proliferation of phylogenomics pipelines designed for infectious disease surveillance and outbreak investigation. Transitioning the use of WGS data out of the research laboratory and into the front lines of surveillance and outbreak response requires user-friendly, reproducible and scalable pipelines that have been well validated. Single Nucleotide Variant Phylogenomics (SNVPhyl) is a bioinformatics pipeline for identifying high-quality single-nucleotide variants (SNVs) and constructing a whole-genome phylogeny from a collection of WGS reads and a reference genome. Individual pipeline components are integrated into the Galaxy bioinformatics framework, enabling data analysis in a user-friendly, reproducible and scalable environment. We show that SNVPhyl can detect SNVs with high sensitivity and specificity, and identify and remove regions of high SNV density (indicative of recombination). SNVPhyl is able to correctly distinguish outbreak from non-outbreak isolates across a range of variant-calling settings, sequencing-coverage thresholds or in the presence of contamination. SNVPhyl is available as a Galaxy workflow, Docker and virtual machine images, and a Unix-based command-line application. SNVPhyl is released under the Apache 2.0 license and available at <http://snvphyl.readthedocs.io/> or at <https://github.com/phac-nml/snvphyl-galaxy>.

DATA SUMMARY

1. Simulated sequence reads used to evaluate SNVPhyl (both for variant identification and contamination) have been deposited in FigShare: <https://doi.org/10.6084/m9.figshare.4294838>.
2. Code used to perform the SNVPhyl evaluations for this study is available on GitHub/Zenodo: <https://github.com/apetkau/snvphyl-validations>; DOI: 10.5281/zenodo.439977.

INTRODUCTION

The high-efficiency and cost-effectiveness of whole-genome sequencing (WGS) using next-generation sequencing technologies is transforming the biomedical landscape. Entire

microbial genomes can be rapidly sequenced and subsequently queried with nucleotide-level resolution, an exciting new ability that far outstrips other traditional microbial typing methods. This powerful new ability has the potential to advance many fields, including in particular the field of infectious disease genomic epidemiology. A number of landmark studies have demonstrated the power of WGS for molecular epidemiology. One notable study is the investigation into the 2010 Haiti cholera outbreak [1–3], where WGS and epidemiological data were used in support of the hypothesis that cholera was introduced to Haiti from United Nations peacekeepers originally infected in Nepal. WGS has supported the investigation of outbreaks of organisms as diverse as *Mycobacterium tuberculosis* [4, 5],

Received 6 February 2017; Accepted 12 April 2017

Author affiliations: ¹National Microbiology Laboratory, Public Health Agency of Canada, Winnipeg, MB R3E 3R2, Canada; ²University of Manitoba, Winnipeg, MB R3T 2N2, Canada; ³Health Canada – Bureau of Microbial Hazards, Ottawa, ON K1A 0K9, Canada; ⁴Lethbridge Research and Development Centre, Lethbridge, AB T1J 4B1, Canada; ⁵Centers for Disease Control and Prevention, Atlanta, GA 30333, USA; ⁶Dalhousie University, Halifax, NS B3H 4R2, Canada; ⁷BC Public Health Microbiology and Reference Laboratory, Vancouver, BC V5Z 4R4, Canada; ⁸Simon Fraser University, Burnaby, BC V5A 1S6, Canada.

*Correspondence: Gary Van Domselaar, gary.vandomselaar@phac-aspc.gc.ca

Keywords: genomic epidemiology; phylogenomics; single nucleotide variation detection; bioinformatics; infectious disease surveillance; bacterial genomics.

Abbreviations: API, application programming interface; hqSNV, high-quality single nucleotide variant; IRIDA, integrated rapid infectious disease analysis; MLST, multilocus sequence typing; NCBI, National Center for Biotechnology Information; SNV, single nucleotide variant; WGS, whole-genome sequencing.

Data statement: All supporting data, code and protocols have been provided within the article or through supplementary data files. Two supplementary figures and three supplementary tables are available with the online Supplementary Material.

Escherichia coli [6] and *Legionella pneumophila* [7]. These high-profile successes have motivated public-health institutions and food-regulatory agencies to incorporate WGS into their routine microbial infectious disease surveillance and outbreak investigation activities. The GenomeTrakr network used by the Centers for Disease Control and the Food and Drug Administration agencies in the USA [8], PulseNet International (<http://www.cdc.gov/pulsenet/next-generation.html>), Statens Serum Institut in Denmark [9], and Public Health England [10] are leading the charge in this area, and have incorporated a variety of analytical approaches to integrate WGS into their infectious disease surveillance activities. Two approaches in particular have emerged as feasible methods for bacterial genomic epidemiology: gene-by-gene methods, which extend the idea of multilocus sequence typing (MLST) to encompass a given organism's entire genome (whole-genome MLST, wgMLST) or core genome (core genome MLST, cgMLST) [11, 12]; and single nucleotide variant (SNV)-based methods (also called single nucleotide polymorphism- or SNP-based methods), which identify variants by comparing a population of target genomes against a reference [13, 14]. Gene-by-gene methods are promising as they are more amenable to assigning consistent sequence types using standardized MLST schemas, but these schemas must be developed, validated, and maintained for each organism. SNV-based methods are popular as they do not require development of MLST schemas, but the variability in SNV-identification methods and reference genome selection means they do not yet produce standard sequence types useful for global communication of circulating infectious disease [12, 14]. Where applicable, these two methods are often combined [15].

A growing number of SNV-based pipelines have been developed (Table 1) and are distributed in the form of web services [16], command-line software [17] or both [14]. Web services provide a user-friendly method of running large-scale analyses, but require the uploading of sequence reads and rely on third-party computing infrastructure, which may be inadequate for the analysis of typically large datasets or due to data privacy concerns. Locally installed pipelines avoid the transfer of large datasets to third-party websites, offer greater control over the execution environment for reproducibility and allow for the incorporation into pre-existing bioinformatics analysis environments. However, locally installed pipelines may require considerable expertise to operate and can have substantial computing requirements. Additionally, for many SNV-based pipelines, recombination detection and removal may require pre-analysis to identify phage and genomic islands in the reference genome, or post-analysis with computationally intensive recombination-detection software such as Gubbins [18] or ClonalFrameML [19] to identify and mask possible recombinant regions. While a large choice of pipelines is available, a systematic comparison of popular SNV pipelines has demonstrated that they generally produce highly concordant phylogenetic trees, but with variation in the particular SNVs identified [20, 21]. However, variation in the installation

IMPACT STATEMENT

The widespread application of whole-genome sequencing (WGS) to investigate infectious disease outbreaks has led to a proliferation of methods for classifying microbial samples based on genomic data. Single-nucleotide variant (SNV)-based methods have been especially useful and a number of SNV phylogenomic pipelines are now available. However, these pipelines often operate under different execution environments, from locally installed command-line applications to web-based applications, and may require multiple analysis software for further processing of the resulting files. We have developed Single Nucleotide Variant Phylogenomics (SNVPhyl) as an SNV-based phylogenomics pipeline that is integrated within the Galaxy platform providing a locally installable environment for phylogenomics analysis within a larger-scale bioinformatics system. We also provide a command-line interface for batch execution and integrate SNVPhyl within the Integrated Rapid Infectious Disease Analysis (IRIDA) genomic epidemiology platform (<http://irida.ca>). We evaluate SNVPhyl's methods for identifying and removing SNVs in recombinant regions and show that SNVPhyl can be used either alone or as input to existing recombination-detection software. We also evaluate SNVPhyl's performance on WGS data from different outbreaks under a number of scenarios, from the presence of low-coverage samples to cross-contamination of closely related samples. We demonstrate the successes and limitations of SNVPhyl under these scenarios, and provide guidance on identifying and handling problematic results.

procedures and execution environments of these pipelines proves challenging for integration into a larger bioinformatics analysis system.

Galaxy [22] is a web-based biological data analysis platform that can be accessed through a publicly available website, a locally installed instance linked to a high-performance compute cluster or a cloud-based environment. Galaxy provides a user-friendly web interface for the construction of data analysis workflows using a mixture of built-in or community developed bioinformatics tools. Additionally, Galaxy provides an API (application programming interface) for automated workflow execution or other automations via external software. These features have encouraged some software developers to integrate Galaxy within larger data analysis systems. Examples of such analysis systems include Integrated Rapid Infectious Disease Analysis (IRIDA; <http://irida.ca>), the Refinery Platform (www.refinery-platform.org/) and the Genomics Virtual Laboratory [23].

The Single Nucleotide Variant Phylogenomics (SNVPhyl) pipeline provides a reference-based SNV discovery and phylogenomic tree-building pipeline along with ancillary tools

Table 1. A comparison of whole-genome phylogenetic software

Name	Input*	Parallel computing†	Distribution‡	Interface§	Reference
CFSAN SNP pipeline	sr	mn, mt	Local	cl	[17]
CSI phylogeny	sr, ag	NA	Web	gui	[16]
kSNP	sr, ag	mt	Local	cl	[44]
Lyve-SET	sr, agr	mn, mt	Local	cl	[21]
NASP	sr, ag	mn, mt	Local	cl	[20]
Parsnp	ag	mt	Local	cl	[45]
PhaME	sr, ag	mt	Local	cl	[46]
REALPHY	sr, ag	mt	Web, local	gui, cl	[14]
Snippy	sr, agr	mt	Local	cl	https://github.com/tseemann/snippy
SNVPhyl	sr	mn, mt	Local	gui, cl	http://snvphyl.readthedocs.io/

*ag, assembled genome; agr, assembled genome supported by generating simulated reads; sr, sequence reads.

†mn, multi-node – provides capability to execute across multiple compute nodes; mt, multi-thread – provides multi-threading capability; NA, not applicable (not locally installable).

‡Local, locally distributed and installable software; web, software provided as a web service.

§cl, command-line interface; gui, graphical user interface.

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

METHODS

SNVPhyl pipeline

The SNVPhyl pipeline (Fig. 1a, b) consists of a set of pre-existing and custom-developed bioinformatics tools for reference mapping, variant discovery and phylogeny construction from identified SNVs. Each stage of the pipeline is implemented as a separate Galaxy tool and the stages are joined together to construct the SNVPhyl workflow. Distribution of the dependency tools for SNVPhyl is managed through the Galaxy Toolshed [24]. Scheduling of each tool is managed by Galaxy, which provides support for execution on a single machine, high-performance computing environments utilizing most major scheduling engines (e.g. Slurm, TORQUE, Open Grid Engine) or cloud-based environments.

Input

SNVPhyl requires as input a collection of microbial WGS datasets, a reference genome and an optional masking file defining regions on the reference genome to exclude from the analysis. Each set of sequencing data consists of either single-end or paired-end reads from an isolate. The reference genome can be a high-quality draft or finished genome, chosen typically to

have high similarity with the collection of genome sequences under analysis. The masking file stores the sequence identifier of the reference genome and the coordinates for any regions where SNVs should be excluded from analysis.

Architecture

Execution of SNVPhyl begins with the 'Repeat Identification' stage. This stage identifies internal repeat regions on the reference genome using MUMmer (v3.23) [25] and generates a masking file containing the locations of repetitive regions to exclude from analysis. This file is concatenated to the user-supplied masking file, if defined, and used in later analysis stages.

The 'Mapping/Variant Calling' stage (detailed in Fig. 1b) aligns the supplied reads to the reference genome using the appropriate mapping mode (paired-end or single-end). Reference mapping is performed using SMALT (version 0.7.5; <https://sourceforge.net/projects/smalt/>), which outputs a read pileup. In the 'Mapping Quality' stage, SNVPhyl evaluates each pileup for the mean coverage across a user-defined proportion of the reference genome (e.g. 10× coverage across at least 80 % of the genome). Any sequenced genomes that do not meet the minimum mean coverage threshold are flagged for further assessment.

The variant calling stages of SNVPhyl use two independent variant callers, FreeBayes (version 0.9.20) [26], and the SAMtools and BCFtools packages [27, 28]. FreeBayes is run using the haploid variant calling mode and the resulting variants are filtered to remove insertions/deletions and split complex variant calls. SAMtools and BCFtools are run independently of FreeBayes and are used to confirm the FreeBayes variant calls and generate base calls for non-variant positions.

The 'Variant Consolidation' stage combines both sets of variant and non-variant (polymorphic and monomorphic) calls into a merged file, flagging mismatches between variant callers. Base calls below the defined minimum read coverage

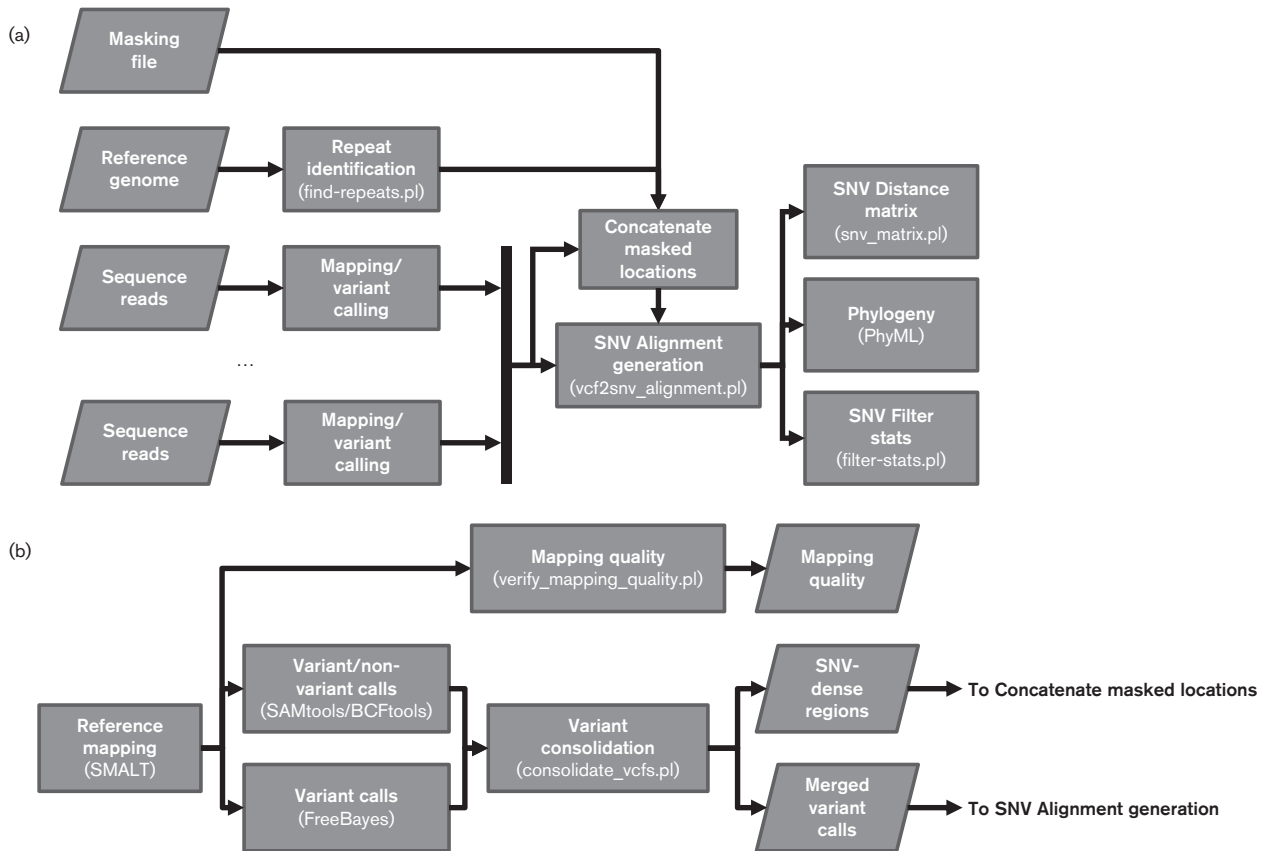


Fig. 1. (a) Overview of the SNVPhyl pipeline. Input to the pipeline is provided as a reference genome, a set of sequence reads for each isolate and an optional list of positions to mask from the final results. Repeat regions are identified on the reference genome and reference mapping followed by variant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

are identified and flagged. The merged base calls are scanned for positions that do not pass the minimum relative SNV abundance (proportion of reads supporting the SNV with respect to the depth of coverage at a site) and minimum mean mapping quality. These base calls are removed from the merged base calls file. The remaining base calls that pass all these criteria are defined as either a high-quality SNV (hqSNV) or a high-quality non-variant base call. The hqSNVs are optionally scanned to identify high-density SNV regions. These regions are identified by passing a sliding window of a given size along the genome and counting the number of SNVs within the window that exceed a given SNV density threshold. The high-density SNV regions are recorded in a tab-delimited file and used to mask potential recombinant regions.

The 'SNV Alignment Generation' stage examines the merged base calls to generate a table of identified variants, and an alignment of hqSNVs and the corresponding high-

quality non-variant bases. The hqSNVs are evaluated and assigned a status using the base calls at the same reference genome position for every isolate. A status of 'valid' is assigned when the base calls from all isolates in the same position pass the minimum criteria (hqSNVs or high-quality non-variants). These base calls are incorporated into the SNV alignment used for phylogeny generation. A status of 'filtered-coverage' is assigned when one or more isolates fail the minimum base coverage threshold at a particular position and the failed isolates' base calls are annotated as '-' (indicating no nucleotide or a gap). A status of 'filtered-mpileup' is assigned when one or more isolates have conflicting base calls between FreeBayes and SAMtools/BCFtools and the conflicting isolates' base calls are annotated as 'N' (indicating any nucleotide non-specifically). A status of 'filtered-invalid' is assigned when the identified hqSNV overlaps one of the masked locations. The hqSNVs, base calls and assigned status are recorded in the SNV table and saved for later inspection. The SNV table can be used to

re-generate the downstream SNV alignment and phylogenetic tree without re-running the computationally intensive reference mapping and variant calling steps.

Output

The final phylogeny is generated using the SNV alignment consisting of hqSNVs with a ‘valid’ status (i.e. the polymorphic positions passing our quality thresholds). This alignment is run through PhyML [29] with the GTR+ γ model as default and tree support values estimated using PhyML’s approximate likelihood ratio test [30]. The SNV alignment is also used to generate an all-against-all SNV distance matrix. This matrix lists the pair-wise SNV distances between every isolate, using only the valid hqSNVs.

Additional files are provided to assist in evaluating the quality of the SNVPhyl analysis. The ‘SNV Filter Stats’ stage summarizes the quality and counts of the identified SNVs. The SNV Alignment Generation stage summarizes the proportion of the reference genome passing all the necessary filters for every isolate – the (non-masked) core genome consisting of both polymorphic and monomorphic positions – as well as the portion of the genome failing any filters or excluded by the masking file.

Simulated data

We evaluated SNVPhyl’s sensitivity and specificity for SNV identification using simulated mutations derived from a reference genome. The closed and finished *E. coli* strain Sakai (NC_002695), along with the two plasmids (NC_002128 and NC_002127), was chosen as the reference genome (combined length of 5 594 477 bp). We constructed a variant genome by randomly mutating 10 000 base locations on the reference genome. We repeated the procedure, using the same 10 000 base locations but different mutations, to generate a total of three variant genomes. We included the unmodified reference genome in the test set to serve as a positive control. The simulated variants for each genome were recorded in a table for later comparisons. The constructed genomes were run through art_illumina (version ChocolateCherryCake) [31] to generate paired-end reads with 2×250 bp length and 30× mean coverage. The resultant reads along with the reference genome were run through SNVPhyl with repeat masking enabled, but with no SNV density filtering.

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

SNV density filtering evaluation

We evaluated SNVPhyl’s ability to mask recombination by comparing the resultant phylogenetic trees and identified SNVs to those detected and removed by the recombination detection software package Gubbins [18]. Our test data consisted of 11 *Streptococcus pneumoniae* genomes along with the reference genome ATCC 700669 (FM211187) that had previously been published [32] and made available as sequence reads by the National Center for Biotechnology Information (NCBI; Table S1, available in the online Supplementary Material) and as a whole-genome alignment constructed via mapping reads to the reference genome (the PMEN1 dataset from <https://sanger-pathogens.github.io/gubbins/>). We downloaded this alignment, appended the reference genome, and processed the resulting file through Gubbins to identify and mask recombinant SNVs. The identified SNVs were filtered to remove gaps and masked recombination (‘-’ and ‘N’ characters) and the resulting SNVs we defined as the ‘truth’ set used to generate the T/F P/N values – defined as for the Simulated Data section. These Gubbins-identified SNVs were also used to construct a phylogenetic tree with PhyML and compared with SNVPhyl’s phylogenetic trees numerically using K tree scores [33] and visually using PhyTools [34]. K tree scores allow for similarity comparisons of many phylogenetic trees against a single reference tree. Each tree is re-scaled by a factor, K, based on the reference tree size and a score is produced taking into account differences in both topology and branch lengths. Comparing the scores of all trees provides a measure of similarity to the reference tree, with more similar trees producing a score closer to 0.

We downloaded sequence reads for the test dataset from NCBI, identifying and combining multiple sequencing runs for each strain to a single set of sequence reads with the help of SRADB [35]. Using the combined sequence reads we ran SNVPhyl under a number of scenarios. For each scenario, we compared the SNVs and phylogenetic trees to the ‘truth’ dataset described above. In the first run, we performed no SNV density filtering. For all subsequent runs, we adjusted the density-filtering parameters to remove SNVs occurring at a density of 2 or more within a moving window of 20, 100, 500, 1000 and 2000 bp. We evaluated an additional scenario using a combination of SNVPhyl and Gubbins for recombination masking. We ran SNVPhyl with no SNV density filtering and incorporated the identified variants into the reference genome to generate a whole-genome alignment consisting of both polymorphic and monomorphic positions, but with non-hqSNVs ignored (i.e. positions with gaps, ambiguous bases or repeats are left as monomorphic). The whole-genome alignment was processed with Gubbins to identify non-recombinant SNVs and to construct a phylogenetic tree.

Parameter optimization

We evaluated SNVPhyl’s parameter settings and resulting accuracy at differentiating outbreak isolates using a set of 59 sequenced and published *Salmonella enterica* serovar

Heidelberg genomes [36], which were previously deposited in the NCBI Sequence Read Archive (Table S2). We chose this dataset as it contained sequence data for strains from several unrelated outbreaks – referred to as ‘outbreak 1’, ‘outbreak 2’ and ‘outbreak 3’ – along with additional background strains, allowing us to evaluate SNVPhyl’s ability to differentiate the outbreak strains under different scenarios. Sequence read data was subsampled with seqtk (<https://github.com/lh3/seqtk>) such that the genome with the least amount of sequence data, SH12-006, was set to 30× mean coverage (calculated as: mean coverage = count of base pairs in all reads/length of reference genome). Other genomes were subsampled to maintain their relative proportion of mean read coverage to SH12-006. *Salmonella* Heidelberg strain SL476 (NC_011083) was selected as the reference genome. We optimized the SNVPhyl parameters for this dataset according to the following four scenarios: (1) adjusting the minimum base coverage parameter used to call a variant, while keeping the number of reads in the dataset fixed; (2) subsampling the reads of a single WGS sample at different mean coverage levels, while keeping the minimum base coverage parameter fixed; (3) adjusting the minimum relative SNV abundance for calling a variant; and (4) adjusting the amount of contamination in the dataset to determine its effect on variant calling accuracy.

In the first scenario, we ran the SNVPhyl pipeline using the default parameters except for the minimum base coverage, which was adjusted to 5×, 10×, 15× and 20×. In the second scenario, we kept the minimum base coverage parameter fixed at 10×, while one of the samples (SH13-001) was subsampled to mean sequencing coverages of 30×, 20×, 15× and 10×. In the third scenario, the minimum relative SNV abundance was adjusted to 0.25, 0.5, 0.75 and 0.9. In the fourth scenario, a sample from outbreak 2 (SH13-001 with mean coverage 71×) was chosen as a candidate for simulating contamination. A sample from the unrelated outbreak 1 (SH12-001) was selected as the source of contaminant reads. The reads were subsampled and combined such that SH13-001 (outbreak 2) remained at 71× mean coverage, but was contaminated with reads from SH12-001 (outbreak 1) at 5, 10, 20 and 30%. All samples were run through SNVPhyl for each of these contamination ratios.

The phylogenetic trees produced by SNVPhyl were evaluated for concordance with the outbreak epidemiological data using the following criteria: (1) all outbreak isolates

group monophyletically; and (2) the SNV distance between any two isolates within an outbreak clade is less than 5 SNVs, a number identified in the previous study [36] as the maximum SNV distance between epidemiologically related samples within these particular outbreaks. Both conditions were tested using the APE package within R [37].

RESULTS

Validation against simulated data

We measured SNVPhyl’s sensitivity and specificity by introducing random mutations along the *E. coli* Sakai reference genome and compared these mutations with those detected by SNVPhyl (Table 2). Of the 10 000 mutated positions introduced, SNVPhyl reported 9116 TPs and 0 FPs resulting in a sensitivity and specificity of 0.91 and 1.0, respectively.

Positions on the reference genome that contain a low-quality base call or exist in repetitive regions are excluded from downstream analysis by SNVPhyl. However, lower-quality variant-containing sites along with variants in repetitive regions are saved by SNVPhyl in the variant table with a ‘filtered’ status. Evaluating the combination of high-quality variants along with the additional low-quality variants recorded by SNVPhyl, we found 457 additional TPs (for a total of 9573) at the expense of 51 FPs, resulting in a sensitivity and specificity of 0.96 and 1.0 (after rounding). Of the 51 FPs, 48 were considered as FPs due to insufficient read coverage in one of the samples to call a high-quality variant; thus, resulting in a call of a gap (‘-’) as opposed to the true base call. Only three of the FPs were a result of miscalled bases with sufficient read coverage, and these occurred in repetitive regions of the genome with high copy numbers (Table S3).

SNV density filtering evaluation

We compared SNVPhyl’s density filtering against the Gubbins software for detection and removal of recombination in a collection of WGS reads from 11 *S. pneumoniae* genomes along with the reference genome ATCC 700669 (Table 3, Fig. S1). We used a previously generated and published whole-genome alignment of these genomes, which we ran through Gubbins to construct a set of 165 non-recombinant SNV-containing positions, which we defined as the TPs used for comparison with SNVPhyl. With no SNV density filtering, SNVPhyl properly identified 142/165 of these SNV-containing sites (TPs), but included 2159 additional

Table 2. SNV simulation results

Comparison	No. of variant columns simulated	No. of non-variant columns	No. of true positives	No. of false positives	No. of true negatives	No. of false negatives	Specificity	Sensitivity
Valid SNVs*	10 000	5 584 477	9116	0	5 575 361	884	1.0	0.91
All SNVs†	10 000	5 584 477	9573	51	5 574 853	427	1.0	0.96

*Valid SNVs – the number of SNV-containing sites detected that passed all thresholds to be considered high quality for every isolate.

†All SNVs – all the SNV-containing sites identified by SNVPhyl, including those where at least one isolate did not have a high-quality base call or sites that were masked by the pipeline.

SNV sites (FPs). These FPs skewed the resulting phylogenetic tree by increasing the length of one of the branches. The phylogenetic tree was compared with the tree produced with Gubbins, resulting in a K tree score of 0.419.

We reanalysed the dataset with high-density SNV masking enabled, using a range of variant density cut-offs. We found the density-filtering criteria of two SNVs in a 500 bp window and two SNVs in a 1000 bp window performed near-equally in producing a phylogenetic tree resembling the tree produced by Gubbins based on the K tree scores of 0.045 and 0.044, both much lower than the score of 0.419 for no SNV density filtering. With these filtering criteria, SNVPhyl identified 133 TPs and 12 FPs (for two SNVs in 500 bp) and 125 TPs and six FPs (for two SNVs in 1000 bp).

We also investigated the effect of generating a whole-genome alignment – by incorporating SNVPhyl-identified variants without SNV density filtering into the reference genome to construct an alignment with both polymorphic and monomorphic positions – for a more thorough analysis with the recombination-detection software Gubbins. We were able to identify 138 TPs in the alignment at the expense of 10 FPs and a K tree score of 0.037, a result closely matching the use SNVPhyl's density filtering criteria.

Parameter optimization

We evaluated SNVPhyl's capability to differentiate between epidemiologically related and unrelated samples using a WGS dataset consisting of 59 *Salmonella enterica* serovar Heidelberg genomes from three unrelated outbreaks. We ran SNVPhyl with this data under a number of scenarios: (1) varying the minimum base coverage required by SNVPhyl to call a variant, (2) subsampling the reads of an individual bacterial sample, (3) varying the minimum relative SNV abundance, and (4) testing the ability to generate accurate phylogenetic trees in the presence of contamination. We tested the SNVPhyl results for phylogenetic concordance to epidemiological data (Table 4, Fig. S2).

For the first scenario, we found that as the minimum base coverage threshold for calling a variant was increased, the

percent of the reference genome identified as part of the core genome and number of SNV-containing sites was reduced (from 95 % core and 317 SNVs to 54 % core and 165 SNVs). At 15× minimum base coverage (81 % core and 262 SNVs) and lower, all three outbreaks grouped into monophyletic clades. Failure occurred at a minimum base coverage of 20× (54 % core and 165 SNVs), where the outbreak 2 isolates failed to constitute a separate clade.

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

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

For the fourth scenario, we examined the robustness of SNVPhyl to cross-contamination of closely related samples. Current methods of contamination detection often focus on taxonomic classification of genomic content [38]. However, contamination by closely related isolates can go undetected, leading to the usage of such contaminated datasets within bioinformatics pipelines. To examine the effect of contamination on SNVPhyl's ability to call hqSNVs, we simulated

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

Case	No. of true positives	No. of false positives	No. of true negatives	No. of false negatives	Sensitivity	Specificity	K tree score
No DF*	142	2159	2 218 849	23	0.861	0.999	0.419
2 in 20†	142	565	2 220 443	23	0.861	1.000	0.425
2 in 100†	142	155	2 220 853	23	0.861	1.000	0.377
2 in 500†	133	12	2 221 005	32	0.806	1.000	0.045
2 in 1000†	125	6	2 221 019	40	0.758	1.000	0.044
2 in 2000†	111	3	2 221 036	54	0.673	1.000	0.063
Gubbins/ SNVPhyl‡	138	10	2 221 002	27	0.836	1.000	0.037

*No DF – a case of no SNV density filtering by SNVPhyl.

†X in Y – masking regions with a density of X variants in Y bases.

‡Gubbins/SNVPhyl – a whole-genome alignment generated from SNVs identified by SNVPhyl and run through Gubbins.

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

No.	Scenario	Parameter/condition	hqSNV	% core*	Differentiated outbreaks
1	Minimum coverage	5×	317	95	Yes
		10×	301	92	Yes
		15×	262	81	Yes
		20×	165	54	No
2	Subsample coverage level	10×†	155	47	No
		15×†	242	76	Yes
		20×†	276	88	Yes
		30×†	299	92	Yes
3	Relative SNV abundance	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

*100 % core = 4 888 768 bp (percentage of reference genome identified as the core genome).

†These represent the mean coverage of one sample after subsampling reads and not the minimum base coverage parameter of SNVPhyl (which is fixed at 10×).

‡100 % contamination represents complete replacement of reads from SH13-001 (at 71× coverage) with SH12-001.

contamination for an isolate in outbreak 2 by an isolate in outbreak 1. We found that SNVPhyl was able to accurately differentiate all three outbreaks with up to 10 % read contamination; however, the number of SNVs dropped from 298 SNVs at 5 % contamination, to 260 SNVs at 20 % contamination, where the failure was due to removal of the majority of unique SNVs that differentiated outbreak 1 from the background isolates.

DISCUSSION

The availability of WGS data from microbial genomes represents a tremendous opportunity for infectious disease surveillance and outbreak response. Emerging analytical methods, such as gene-by-gene or SNV-based methods, require that bioinformatics pipelines be designed with usability by non-bioinformaticians in mind and that can be easily incorporated into existing systems. An overview of current phylogenomic methods appears in [39] and a comparison of SNVPhyl's design with that of other popular pipelines appears in Table 1. A detailed investigation comparing the performance of SNVPhyl with other pipelines is the subject of a separate paper [21]. We designed SNVPhyl to be both flexible and scalable in its usage in order to meet the needs and capabilities of most laboratories. SNVPhyl gains much of this flexibility through its implementation as a Galaxy workflow, which enables execution in environments from single machines to high-scale computer

clusters, from third-party web-based environments to local installations. Galaxy provides a user-friendly interface but also provides an API, which is used to implement a command-line interface for SNVPhyl. The SNVPhyl pipeline is also integrated within the IRIDA platform (<http://irida.ca>), which provides an integrated 'push-button' system for genomic epidemiology. However, implementing SNVPhyl through Galaxy has some disadvantages. Notably, Galaxy is more complex and, thus, more cumbersome to install than a simpler command-line-based pipeline. To address this, we have made SNVPhyl available as simple to install virtual machine and Docker images, although these options are not straightforward to implement in a high-performance computing environment.

Several factors can influence the ability to accurately call SNVs when using a reference mapping approach [40]. In addition, there are aspects of the datasets – such as recombination and population diversity – that can influence the phylogenetic analysis of identified SNVs. To assist in selecting proper parameters for SNVPhyl and gauging performance on different datasets, we have assessed SNVPhyl under a variety of situations: SNV calling accuracy with simulated data, recombination masking, and the ability to differentiate outbreak isolates from non-outbreak isolates under differing parameters and data qualities.

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

Phylogenetics assumes descent with modification, but recombination violates this assumption and its presence can confound the resulting phylogeny leading to misinterpretations on the clonal relationship of isolates [41]. Recombination detection software exists and can be used to account for recombination during the construction of phylogenetic trees [18, 19, 42]. These programs are most effective for the detection of recombination in closely related organisms, such as a collection of bacteria in an epidemiological investigation. However, they require the pre-construction of whole-genome alignments and can only be run on a single machine, which limits their utility for routine application to large collections of WGS reads.

SNVPhyl implements a basic but rapid method for detection and masking of recombinant sites by searching for SNV-dense regions above a defined density in a sliding

window. We evaluated SNVPhyl's recombination-masking method in comparison to the Gubbins software package, which was run on a previously generated whole-genome alignment (Table 3, Fig. S1). We found that SNVPhyl removes the majority of recombinant SNVs (from 2159 SNVs with no recombination masking to six SNVs when masking regions with two SNVs in a 1000 bp window). However, SNVPhyl also removes some non-recombinant SNVs (reduced from 142 SNVs with no masking to 125 SNVs with two SNVs in a 1000 bp window). Removal of a greater number of recombinant SNVs is possible by increasing the window size, but this removes additional non-recombinant SNVs and reduces the information available in the phylogenetic tree and so concordance with other recombination-masking procedures (based on K tree scores).

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

A proper interpretation of the produced phylogenetic trees and SNV distances for associating closely related isolates requires knowledge of when to trust the results and when additional parameter or data adjustments are necessary. To assist in defining these criteria, we evaluated the performance of SNVPhyl at clearly delineating different outbreak clades across four different scenarios (Table 4, Fig. S2a–d).

In both the first and second scenarios, we examined the effect of sequencing coverage on identifying enough SNVs to properly differentiate outbreak isolates. In the first scenario, we adjusted the minimum base coverage required to call a SNV from $5\times$ to $20\times$ without any additional subsampling of reads. We found that SNVPhyl succeeded in differentiating outbreak isolates at coverages up to $15\times$, but at a minimum base coverage of $20\times$ SNVPhyl failed to differentiate the outbreak isolates due to removal of too many SNVs (from 317 SNVs to 165 SNVs). In the second scenario, we subsampled one of the isolates along the mean read coverage values from $30\times$ to $10\times$, while keeping the minimum base coverage parameter in SNVPhyl fixed at $10\times$. We

found SNVPhyl succeeded in differentiating outbreak isolates at a mean coverage of $15\times$ and above, but failed to differentiate outbreak isolates at a mean coverage of $10\times$ due to removal of too many SNVs (reduced from 299 SNVs to 155 SNVs). Both cases show that a high base coverage threshold for variant calling relative to the mean coverage of the lowest sample leads to falsely identifying samples as being related due to removal of too many SNVs ($20\times$ minimum base coverage/ $30\times$ lowest sample mean coverage for failure in the first scenario, and $10\times$ minimum base coverage/ $10\times$ lowest sample mean coverage for failure in the second scenario). However, a high minimum base coverage threshold or too little sequencing data can be detected by examining the percentage of the reference genome considered as part of the core genome by SNVPhyl. A low value can indicate either a poorly related reference genome or that large portions of the genomes are removed from the analysis (a drop from 95 to 54 % in the first scenario and 92 to 47 % in the second scenario). We would recommend searching for such low values in the percent core to gauge whether or not base coverage (or possibly reference genome selection) is an issue for the SNVPhyl results.

In the third scenario (Table 4, Fig. S2c), we adjusted the relative SNV abundance among values from 0.25 to 0.9. We found that SNVPhyl successfully differentiated outbreak isolates above a proportion of 0.5, but at a proportion of 0.5 the maximum SNV distance between isolates within an outbreak exceeded our threshold of less than 5 SNVs. However, unlike the minimum base coverage value, the percent of the reference genome identified as the core genome remained the same (92 %). We recommend keeping this setting fixed at a higher value, with the default set at 0.75.

In the fourth scenario, we simulated contamination between two closely related isolates from two different outbreaks by mixing reads at differing proportions (Table 4, Fig. S2d). Our findings indicate that SNVPhyl is able to handle low amounts of mixed sample contamination (up to 10 %). A higher proportion of contaminated reads can lead to removal of SNVs due to not meeting quality thresholds (from 298 SNVs with 5 % contamination to 260 SNVs at 20 % contamination where failure occurred) and so incorrectly inferring relatedness between samples. Similar to the third scenario, the percentage of the reference genome identified as the core genome remained fixed at 92 %. While SNVPhyl is able to differentiate outbreak isolates at low levels of contamination, SNVPhyl cannot be used to evaluate the degree of contamination. Thus, we would not recommend the straightforward application of SNVPhyl to contaminated datasets without further assessment of the degree of contamination, either through taxonomic identification software such as Kraken [43] or, for closely related isolates, through inspection of the variant calling and read pileup information provided by SNVPhyl.

Our analysis suggests that great care must be taken to reduce sources of noise in genome-wide SNV analysis. Some of this noise relates to quality thresholds for calling

hqSNVs, of which a careful balance is required to eliminate FPs without removal of too many true variants. Other sources include aspects of the WGS datasets or organisms under study such as the presence of contamination or recombination. The studied cases highlight how SNVPhyl is able to produce phylogenetic trees consistent with existing software and epidemiological data under a wide variety of data qualities, and demonstrate when to be sceptical of the results based on additional information generated by SNVPhyl.

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

Funding information

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

Acknowledgements

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

Conflicts of interest

The authors declare that there are no conflicts of interest.

Data bibliography

1. Bekal S, Berry C, Reimer AR, Van Domselaar G, Beaudry G et al. Sequence Read Archive PRJNA305824 (2015).
2. Bergholz TM, Wick LM, Qi W, Riordan JT, Ouellette LM et al. GenBank NC_002695.1 (2007).
3. Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J et al. European Nucleotide Archive ERR016678, ERR016679, ERR016671, ERR019725, ERR019714, ERR016851, ERR019721, ERR016858, ERR019732, ERR016859, ERR019733, ERR019722, ERR019734, ERR019723, ERR016860, ERR019715, ERR019726, ERR016852, ERR016681, ERR016720, ERR016721 (2011).
4. Croucher NJ, Walker D, Romero P, Lennard N, Paterson GK et al. GenBank FM211187.1 (2009).
5. Fricke WF, Mammel MK, McDermott PF, Tartera C, White DG et al. GenBank NC_011083.1 (2011).
6. Makino K, Ishii K, Yasunaga T, Hattori M, Yokoyama K et al. GenBank NC_002128.1 (1998).
7. Makino K, Ishii K, Yasunaga T, Hattori M, Yokoyama K et al. GenBank NC_002127.1 (1998).
8. Petkau A, Mabon P, Sieffert C, Knox N, Cabral J et al. FigShare <http://dx.doi.org/10.6084/m9.figshare.4294838> (2016).
9. Petkau A, Mabon P, Sieffert C, Knox N, Cabral J et al. GitHub/Zenodo <https://doi.org/10.5281/zenodo.439977> (2017).

References

1. Hendriksen RS, Price LB, Schupp JM, Gillece JD, Kaas RS et al. Population genetics of *Vibrio cholerae* from Nepal in 2010: evidence on the origin of the Haitian outbreak. *MBio* 2011;2:e00157-11.
2. Katz LS, Petkau A, Beaulaurier J, Tyler S, Antonova ES et al. Evolutionary dynamics of *Vibrio cholerae* O1 following a single-source introduction to Haiti. *MBio* 2013;4:e00398-13.
3. Frerichs RR, Keim PS, Barraix R, Piarroux R. Nepalese origin of cholera epidemic in Haiti. *Clin Microbiol Infect* 2012;18:E158-E163.
4. Gardy JL, Johnston JC, Ho Sui SJ, Cook VJ, Shah L et al. Whole-genome sequencing and social-network analysis of a tuberculosis outbreak. *N Engl J Med* 2011;364:730-739.
5. Roetzer A, Diel R, Kohl TA, Rückert C, Nübel U et al. Whole genome sequencing versus traditional genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a longitudinal molecular epidemiological study. *PLoS Med* 2013;10:e1001387.
6. Holmes A, Allison L, Ward M, Dallman TJ, Clark R et al. Utility of whole-genome sequencing of *Escherichia coli* O157 for outbreak detection and epidemiological surveillance. *J Clin Microbiol* 2015; 53:3565-3573.
7. Sánchez-Busó L, Comas I, Jorques G, González-Candelas F. Recombination drives genome evolution in outbreak-related *Legionella pneumophila* isolates. *Nat Genet* 2014;46:1205-1211.
8. Allard MW, Strain E, Melka D, Bunning K, Musser SM et al. Practical value of food pathogen traceability through building a whole-genome sequencing network and database. *J Clin Microbiol* 2016; 54:1975-1983.
9. Franz E, Gras LM, Dallman T. Significance of whole genome sequencing for surveillance, source attribution and microbial risk assessment of foodborne pathogens. *Curr Opin Food Sci* 2016;8: 74-79.
10. Ashton PM, Nair S, Peters TM, Bale JA, Powell DG et al. Identification of *Salmonella* for public health surveillance using whole genome sequencing. *PeerJ* 2016;4:e1752.
11. Maiden MC, Jansen van Rensburg MJ, Bray JE, Earle SG, Ford SA et al. MLST revisited: the gene-by-gene approach to bacterial genomics. *Nat Rev Microbiol* 2013;11:728-736.
12. Moura A, Criscuolo A, Pouseele H, Maury MM, Leclercq A et al. Whole genome-based population biology and epidemiological surveillance of *Listeria monocytogenes*. *Nat Microbiol* 2016;2:16185.
13. Kwong JC, Mercouliou K, Tomita T, Easton M, Li HY et al. Prospective whole-genome sequencing enhances national surveillance of *Listeria monocytogenes*. *J Clin Microbiol* 2016;54:333-342.
14. Bertels F, Silander OK, Pachkov M, Rainey PB, van Nimwegen E. Automated reconstruction of whole-genome phylogenies from short-sequence reads. *Mol Biol Evol* 2014;31:1077-1088.
15. Jackson BR, Tarr C, Strain E, Jackson KA, Conrad A et al. Implementation of nationwide real-time whole-genome sequencing to enhance listeriosis outbreak detection and investigation. *Clin Infect Dis* 2016;63:380-386.
16. Kaas RS, Leekitcharoenphon P, Aarestrup FM, Lund O. Solving the problem of comparing whole bacterial genomes across different sequencing platforms. *PLoS One* 2014;9:e104984.
17. Davis S, Pettengill JB, Luo Y, Payne J, Shpuntoff A et al. CFSAN SNP Pipeline: an automated method for constructing SNP matrices from next-generation sequence data. *PeerJ Comput Sci* 2015; 1:e20.
18. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA et al. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res* 2015;43:e15.
19. Didelot X, Wilson DJ. ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLoS Comput Biol* 2015;11:e1004041.
20. Sahl JW, Lemmer D, Travis J, Schupp JM, Gillece JD et al. NASP: an accurate, rapid method for the identification of SNPs in WGS

- datasets that supports flexible input and output formats. *Microb Genom* 2016;2:e000074.
21. Katz LS, Griswold T, Williams-Newkirk AJ, Wagner D, Petkau A et al. A comparative analysis of the Lyve-SET phylogenomics pipeline for genomic epidemiology of foodborne pathogens. *Front Microbiol* 2017;8:375.
 22. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D et al. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res* 2016;44:W3–W10.
 23. Afgan E, Sloggett C, Goonasekera N, Makunin I, Benson D et al. Genomics virtual laboratory: a practical bioinformatics workbench for the cloud. *PLoS One* 2015;10:e0140829.
 24. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E et al. Dissemination of scientific software with Galaxy ToolShed. *Genome Biol* 2014;15:403.
 25. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M et al. Versatile and open software for comparing large genomes. *Genome Biol* 2004;5:R12.
 26. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. *arXiv* 2012; arXiv:1207.3907
 27. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J et al. The sequence alignment/map format and SAMtools. *Bioinformatics* 2009;25:2078–2079.
 28. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* 2011;27:2987–2993.
 29. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst Biol* 2003;52:696–704.
 30. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W et al. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* 2010;59:307–321.
 31. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator. *Bioinformatics* 2012;28:593–594.
 32. Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J et al. Rapid pneumococcal evolution in response to clinical interventions. *Science* 2011;331:430–434.
 33. Soria-Carrasco V, Talavera G, Igea J, Castresana J. The K tree score: quantification of differences in the relative branch length and topology of phylogenetic trees. *Bioinformatics* 2007;23:2954–2956.
 34. Revell LJ. Phytools: an R package for phylogenetic comparative biology (and other things). *Methods Ecol Evol* 2012;3:217–223.
 35. Zhu Y, Stephens RM, Meltzer PS, Davis SR. SRADB: query and use public next-generation sequencing data from within R. *BMC Bioinformatics* 2013;14:19.
 36. Bekal S, Berry C, Reimer AR, van Domselaar G, Beaudry G et al. Usefulness of high-quality core genome single-nucleotide variant analysis for subtyping the highly clonal and the most prevalent *Salmonella enterica* serovar Heidelberg clone in the context of outbreak investigations. *J Clin Microbiol* 2016;54:289–295.
 37. Paradis E, Claude J, Strimmer K. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics* 2004;20:289–290.
 38. Koren S, Treangen TJ, Hill CM, Pop M, Phillippy AM. Automated ensemble assembly and validation of microbial genomes. *BMC Bioinformatics* 2014;15:126.
 39. Lynch T, Petkau A, Knox N, Graham M, van Domselaar G. A primer on infectious disease bacterial genomics. *Clin Microbiol Rev* 2016;29:881–913.
 40. Olson ND, Lund SP, Colman RE, Foster JT, Sahl JW et al. Best practices for evaluating single nucleotide variant calling methods for microbial genomics. *Front Genet* 2015;6:235.
 41. Croucher NJ, Harris SR, Grad YH, Hanage WP. Bacterial genomes in epidemiology—present and future. *Philos Trans R Soc Lond B Biol Sci* 2013;368:20120202.
 42. Marttinen P, Hanage WP, Croucher NJ, Connor TR, Harris SR et al. Detection of recombination events in bacterial genomes from large population samples. *Nucleic Acids Res* 2012;40:e6.
 43. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol* 2014;15:R46.
 44. Gardner SN, Slezak T, Hall BG. kSNP3.0: SNP detection and phylogenetic analysis of genomes without genome alignment or reference genome. *Bioinformatics* 2015;31:2877–2878.
 45. Treangen TJ, Ondov BD, Koren S, Phillippy AM. The Harvest suite for rapid core-genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome Biol* 2014;15:524.
 46. Ahmed SA, Lo C, Li P, Davenport KW, Chain PSG et al. From raw reads to trees: whole genome SNP phylogenetics across the tree of life. *bioRxiv* 2015; doi:10.1101/032250.

Five reasons to publish your next article with a Microbiology Society journal

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

Find out more and submit your article at microbiologyresearch.org.