1   **Title:** The Northern Arizona SNP Pipeline (NASP): accurate, flexible, and rapid

2   identification of SNPs in WGS datasets

3

4   Jason W. Sahl[1,2]*, Darrin Lemmer[2]*, Jason Travis[2], James M. Schupp[2], John D.

5   Gillece[2], Maliha Aziz[2,3], Elizabeth M. Driebe[2], Kevin Drees[1,4], Nathan Hicks[2],

6   Charles H.D. Williamson[1], Crystal Hepp[1], David Smith[2], Chandler Roe[2], David M.

7   Engelthaler[2], David M. Wagner[1], Paul Keim[1,2]

8

9   [1]Center for Microbial Genetics and Genomics, Northern Arizona University,

10  Flagstaff, AZ; [2]Translational Genomics Research Institute, Flagstaff, AZ; [3]The

11  George Washington University, Washington D.C.; [4]University of New Hampshire

12

13  *Corresponding authors, Equal contribution to the manuscript

14

15  <u>**Abstract**</u>

16      Whole genome sequencing (WGS) of bacteria is becoming standard practice

17  in many laboratories. Applications for WGS analysis include phylogeography and

18  molecular epidemiology, using single nucleotide polymorphisms (SNPs) as the

19  unit of evolution. The Northern Arizona SNP Pipeline (NASP) was developed as

20  a reproducible pipeline that scales well with the large amount of WGS data

21  typically used in comparative genomics applications. In this study, we

22  demonstrate how NASP compares to other tools in the analysis of two real

23  bacterial genomics datasets and one simulated dataset. Our results demonstrate

24  that NASP produces comparable, and often better, results to other pipelines, but

25  is much more flexible in terms of data input types, job management systems,

26  diversity of supported tools, and output formats. We also demonstrate differences

27  in results based on the choice of the reference genome and choice of inferring

28  phylogenies from concatenated SNPs or alignments including monomorphic

29  positions. NASP represents a source-available, version-controlled, unit-tested

30  method and can be obtained from tgennorth.github.io/NASP.

31

## Introduction

Whole genome sequence (WGS) data from bacteria are rapidly increasing in public databases and have been used for outbreak investigations [1, 2], associating phylogeny with serology [3], as well as phylogeography [4]. WGS data are frequently used for variant identification, especially with regards to single nucleotide polymorphisms (SNPs). SNPs are used because they provide stable markers of evolutionary change between genomes [5]. Accurate and reliable SNP identification requires the implementation of methods to call, filter, and merge SNPs with tools that are version controlled, unit tested, and validated [6].

Multiple pipelines are currently available for the identification of SNPs from diverse WGS datasets, although the types of supported input files differ substantially. There are few pipelines that support the analysis of both raw sequence reads as well as genome assemblies. The ISG pipeline [7] calls SNPs from both raw reads, primarily from the Illumina platform, and genome assemblies, but wasn't optimized for job management systems and only exports polymorphic positions. While only polymorphic positions may be adequate for many studies, including monomorphic positions in the alignment has been shown to be important for various phylogenetic methods. A commonly used SNP analysis software method is kSNP, which has been discussed in three separate publications [8-10]. The most recent version of kSNP (v3) doesn't directly support the use of raw reads in the identification of SNPs. kSNP is a reference-independent approach in which all kmers of a defined length are compared to identify SNPs. The all-versus-all nature of the algorithm can result in a large RAM footprint and can stall on hundreds of bacterial genomes [7]. Finally, REALPHY was published as a method to identify SNPs using multiple references and then merging the results [11]. The authors claim that single reference based methods bias the results, especially from mapping raw reads against a divergent reference genome.

Additional methods have also been published that only support specific input formats. Parsnp is a method that can rapidly identify SNPs from the core

2

63  genome, but currently only processes closely related genome assemblies [12].

64  SPANDx is a method that only supports raw reads, but does run on a variety of

65  job management systems [13]. The program lyve-SET has been used in

66  outbreak investigations and uses raw or simulated reads to identify SNPs [14].

67  Finally, the CFSAN SNP-pipeline is a published pipeline from the United States

68  Food and Drug Administration that only supports the use of raw reads [15]. There

69  have been no published comparative studies to compare the functionality of

70  these pipelines on a range of test datasets.

71  In this study, we describe the Northern Arizona SNP Pipeline (NASP). NASP

72  is a source-available, unit-tested, version-controlled method to rapidly identify

73  SNPs that works on a range of job management systems, incorporates multiple

74  read aligners and SNP callers, works on both raw reads and genomes

75  assemblies, calls both monomorphic and polymorphic positions, and has been

76  validated on a range of diverse datasets. We compare NASP with other methods,

77  both reference-dependent and reference-independent, in the analysis of three

78  reference datasets.

79

80  **Methods**

81

82  NASP is implemented in Python and Go. NASP accepts multiple file formats as

83  input, including ".fasta", ".sam", ".bam", ".vcf", ".fastq", and "fastq.gz". NASP can

84  either function through a question/answer command line interface designed for

85  ease of use, or through an argument-driven command-line interface. NASP was

86  developed to work on job management systems including Torque, Slurm, and

87  Sun/Oracle Grid Engine (SGE); a single node solution is available for NASP as

88  well, but is not optimal. If filtering of duplicate regions in the reference genome is

89  requested, the reference is aligned against itself with NUCmer [16]. These

90  duplicated regions are then masked from downstream analyses, although are still

91  available for investigation. If external genome assemblies are supplied, they are

92  also aligned against the reference genome with NUCmer and SNPs are identified

93  by a direct one-to-one mapping of the query to the reference. In the case of

94    duplications in the query but not the reference, all copies are aligned and any

95    differences at any given base are masked with an "N" character to identify it as

96    ambiguous.

97    If raw reads are supplied, they can be adapter and/or quality trimmed with

98    Trimmomatic [17]. Raw or trimmed reads are aligned against a FASTA-formatted

99    reference using one of the supported short read aligners, including BWA-MEM

100   [18], Novoalign, bowtie2 [19] and SNAP [20]. A binary alignment map (BAM) file

101   is created with Samtools [21] and SNPs can be called with multiple SNP callers,

102   including the UnifiedGenotyper method in GATK [22, 23], Samtools, SolSNP

103   (http://sourceforge.net/projects/solsnp/), and VarScan [24]. Positions that fail a

104   user-defined depth and proportion threshold are filtered from downstream

105   analyses but are retained in the "master" matrices. A workflow of the NASP

106   pipeline is shown in Figure 1 and a summary is shown in Supplemental Table 1.

107   The results of the pipeline can include up to four separate SNP matrices. The

108   first matrix is the master matrix (master.tsv), which includes all calls, both

109   monomorphic and polymorphic, across all positions in the reference with no

110   positions filtered or masked; positions that fall within duplicated regions are

111   shown in this matrix, although they are flagged as duplicated. An optional second

112   matrix (master_masked.tsv) can also be produced. This matrix is the same as

113   the master matrix, although any position that fails a given filter (minimum depth,

114   minimum proportion) is masked with an "N", whereas calls that could not be

115   made are given an "X"; this matrix could be useful for applications where all high-

116   quality, un-ambiguous positions should be considered. The third matrix

117   (missingdata.tsv) includes only positions that are polymorphic across the sample

118   set, but can include positions that are missing in a subset of genomes and not

119   found in duplicated regions; these SNPs have also been processed with the

120   minimum depth and proportion filters and are still high quality calls. The last

121   matrix (bestsnp.tsv) is a matrix with only polymorphic, non-duplicated, clean calls

122   (A,T,C,G) that pass all filters across all genomes. FASTA files are automatically

123   produced that correspond to the bestsnp and missingdata matrices.

124      In addition to the matrices and FASTA files, NASP produces statistics that

125      can be useful for the identification of potentially problematic genomes, such as

126      low coverage or mixtures of multiple strains. These statistics can also be used for

127      determining the size of the core genome, including both monomorphic and

128      polymorphic positions, of a given set of genomes.

129      Post matrix scripts are included with NASP in order to convert between file

130      formats, remove genomes and/or SNPs, provide functional SNP information, and

131      to convert into formats that can be directly accepted by other tools, such as Plink

132      [25], a method to conduct genome wide association studies (GWAS).

133      Documentation for all scripts is included in the software repository.

134

135      **Test datasets**. To demonstrate the speed and functionality of the NASP pipeline,

136      three datasets were selected. The first includes a set of 21 *Escherichia coli*

137      genome assemblies used in other comparative studies [11, 26] (Supplemental

138      Table 2). REALPHY was run on self-generated single-ended simulated reads,

139      100bp in length. Additional pipelines were run with paired-end reads generated

140      by ART chocolate cherry cake [27], using the following parameters: -l 100 -f 20 -p

141      -ss HS25 -m 300 -s 50. Unless otherwise noted, the reference genome for SNP

142      comparisons was K-12 MG1655 (NC_000913) [28]. All computations were

143      performed on a single node, 16-core server with 48Gb of available RAM. For

144      kSNP, the optimum k value was selected by the KChooser script included with

145      the repository.

146      The second dataset includes a set of 15 *Yersinia pestis* genomes from North

147      America (Supplemental Table 3). For those external SNP pipelines that only

148      support raw reads, simulated reads were generated from genome assemblies

149      with ART. A set of SNPs (Supplemental Table 4) has previously been

150      characterized on these genomes with wet-bench methods (unpublished). This set

151      was chosen to determine how many verified SNPs could be identified by different

152      SNP pipelines. All computations were performed on a single node, 16-core

153      server with 48Gb of available RAM.

154  The last dataset includes simulated data from *Yersinia pestis*. Reads and

155  assemblies from 133 *Y. pestis* genomes [29] were downloaded from public

156  databases and processed with NASP using the CO92 genome as the reference

157  to produce a reference phylogeny for WGS data simulation. Assemblies and

158  reads were simulated from this reference phylogeny and a reference genome

159  (CO92   chromosome)   using   TreeToReads

160  (https://github.com/snacktavish/TreeToReads), introducing 3501 mutations. A

161  phylogeny was inferred from the concatenated SNP alignment (3501 simulated

162  SNPs produced by TreeToReads) with RAxML v8 to provide a 'true' phylogeny

163  for the simulated data. Simulated reads (250bp) and assemblies were both

164  processed with pipelines to identify how many of these introduced SNPs could be

165  identified.

166  To test the scalability of NASP on genome assemblies, a set of 3520 *E. coli*

167  genomes was selected (Supplemental Table 5). Genomes were randomly

168  selected   with   a   python   script

169  (https://gist.github.com/jasonsahl/990d2c56c23bb5c2909d) at various levels

170  (100-1000) and processed with NASP. In this case, NASP was run on multiple

171  nodes across a 31-node cluster at Northern Arizona University. The elapsed time

172  was reported only for the step where aligned files are compiled into the resulting

173  matrix. Time required for the other processes is dependent on the input file type

174  and the amount of available resources on a HPC cluster.

175

176  **External SNP pipelines**. Multiple SNP pipelines, both reference-dependent and

177  reference-independent, were compared with NASP, including kSNP v3.9.1 [10],

178  ISG v0.16.10-3 [7], Parsnp v1.2 [12], REALPHY v112 [11], SPANDx v2.7 [13],

179  Mugsy v1r2.2 [30], lyve-SET v1.1.6 [31], and CFSAN (https://github.com/CFSAN-

180  Biostatistics/snp-pipeline). Exact commands used to run each method are shown

181  in Supplemental Data File 1. An overview of all tested methods is shown in Table

182  1. Most of the methods output FASTA files, which were used to infer

183  phylogenies. For Mugsy, the MAF file was converted to FASTA with methods

184  described previously [32].

185

186 **Phylogenetics**. Phylogenies were inferred using a maximum likelihood algorithm
187 implemented in RAxML v8.1.7 [33], except where noted. The exact commands
188 used to infer the phylogenies are shown in Supplemental Data File 1. Tree
189 topologies were also compared on the same input data. Commands to infer
190 these phylogenies using FastTree2 [34], ExaBayes [35], and Parsimonator
191 (github.com/stamatak/Parsimonator-1.0.2) are shown in Supplemental Data File
192 1.

193

194 **Dendrogram of multiple methods**. To visually represent how well different
195 methods relate, a dendrogram was generated. Each phylogeny was compared
196 against a maximum likelihood phylogeny inferred from the reference test set with
197 compare2trees. A UPGMA dendrogram was then calculated with Phylip [36] on
198 the resulting similarity matrix.

199

200 **Results**

201

202 **Pipeline functionality and post-matrix scripts**. NASP is a reference-
203 dependent pipeline that can incorporate both raw reads and assemblies in the
204 SNP discovery process; NASP was not developed for the identification and
205 annotation of short insertions/deletions (indels). NASP can use multiple aligners
206 and SNP callers to identify SNPs and the consensus calls can be calculated
207 across all methods. An additional strength of NASP is that it can run on multiple
208 job management systems as well as on a single node. A complete workflow of
209 the NASP method is shown in Figure 1. Several post-matrix scripts are included
210 with NASP in order to convert between file formats, including generating input
211 files for downstream pipelines, including Plink [25]. An additional script can
212 annotate a NASP SNP matrix using SnpEff [37] to provide functional information
213 for each SNP.

214

215 **NASP run time scalability**. To visualize how NASP scales on processing

216 genome assemblies, a set of 3520 *E. coli* genomes was sampled at 100 genome

217 intervals and processed with NASP with 10 replicates. The results demonstrate

218 that the matrix building step in NASP scales linearly with the processing of

219 additional genomes (Figure 2A). The memory footprint of this step also scales

220 linearly (Figure 2B) and doesn't exceed 4Gb on a large set of genomes (n=1000)**.**

221 If raw reads are used, additional time is required for the alignment and SNP

222 calling methods, and the overall wall time would scale with the number of reads

223 that needed to be processed. The matrix-building step, where assemblies and

224 VCF files are merged into the matrix, would scale linearly regardless of the SNP

225 identification method chosen.

226

227 **Pipeline comparisons on *E. coli* genomes data set.** To test differences

228 between multiple pipelines, a set of 21 *E. coli* genomes used in other

229 comparative genomics studies [11, 26] were downloaded and processed with

230 Parsnp, SPANDx, kSNPv3, ISG, REALPHY, CFSAN, lyve-SET, Mugsy, and

231 NASP. For methods that do not support genome assemblies, paired end reads

232 were simulated with ART, while single end reads were used by REALPHY, as

233 this method is integrated into the pipeline.

234 To identify how well the simulated paired end reads represent the finished

235 genomes, a NASP run was conducted on a combination of completed genome

236 assemblies as well as simulated raw reads. The phylogeny demonstrates that

237 assemblies and raw reads fall into identical locations (Supplemental Figure 1),

238 suggesting that the paired end reads are representative of the finished genome

239 assemblies.

240 The authors of REALPHY assert that their analysis of this dataset

241 demonstrates the utility of using their approach to avoid biases in the use of a

242 single reference genome by using multiple references [11]. However, in our tests,

243 we could only get REALPHY to complete when using a single reference. To test

244 differences between methods, SNPs were identified with multiple reference-

245 dependent and -independent methods, and maximum likelihood (ML)

8

246 phylogenies were compared. The results demonstrate that all methods, with the

247 exception of kSNPv3 and lyve-SET, returned a phylogeny with the same

248 topology as the published phylogeny [11] (compare2trees topological score =

249 100%) (Table 2). The run wall time demonstrates that most other methods were

250 significantly faster than REALPHY (Table 2), even when REALPHY was called

251 against a single reference. Wall time comparisons between methods are

252 somewhat problematic, as some pipelines infer phylogenies and others, including

253 NASP, do not. Additionally, using raw reads is generally expected to be slower

254 than using a draft or finished genome assembly. Finally, some methods are

255 optimized for job management systems, whereas others were designed to run on

256 a single node. For these comparisons, all methods that have single node support

257 were run on a single node. Only SPANDx seems to be dependent on job

258 management systems and could not be successfully run on a single node.

259 One of the other assertions of the REALPHY authors is that phylogenies

260 reconstructed using an alignment of concatenated SNPs are unreliable [11, 26],

261 especially with regards to branch length biases [38]. However, the phylogeny

262 inferred from a NASP alignment of monomorphic and polymorphic sites was in

263 complete agreement with the topology of the phylogeny inferred from a

264 concatenation of only SNPs (compare2trees topological score = 100%); tree

265 lengths were indeed variable by using these two different input types using the

266 same substitution model (Supplemental Figure 2). We also employed an

267 ascertainment bias correction (Lewis correction) implemented in RaxML [38], in

268 order to correct for the use of only polymorphic sites, and found no difference

269 between tree topologies using substitution models that did not employ this

270 correction (data not shown). For this dataset of genomic *E. coli* assemblies, there

271 appears to be no effect of using a concatenation of polymorphic sites on the

272 resulting tree topology, although branch lengths were affected compared to an

273 alignment containing monomorphic sites.

274 To understand how the choice of the reference affects the analysis, NASP

275 was also run using *E. coli* genome assemblies and simulated reads against the

276 outgroup, *E. fergusonii*, as the reference. The results demonstrate that the same

277   tree topology was obtained by using a different, and much more divergent,

278   reference (compare2trees topology score = 100%). However, in both cases,

279   fewer SNPs were identified by using a divergent reference (Table 2).

280       Some authors suggest that reference-independent approaches are less

281   biased and more reliable than reference dependent-approaches [8]. For the case

282   of this *E. coli* dataset, the phylogeny inferred by Mugsy, a reference-independent

283   approach, was in topological agreement with other reference-dependent

284   approaches (Table 2). In fact, kSNPv3 was one of the only methods that returned

285   a tree phylogeny that was inconsistent with all other methods (Table 2); an

286   inconsistent kSNP phylogeny has also been reported in the analysis of other

287   datasets [15]. To analyze this further, we identified SNPs (n=826) from the NASP

288   run using simulated paired-end reads that were uniquely shared on a branch of

289   the phylogeny that defines a monophyletic lineage (Supplemental Figure 3). We

290   then calculated how many of these SNPs were identified by all methods and

291   found widely variable results (Table 2). Using kSNP with only core genome SNPs

292   identified only 5 of these SNPs, which explains the differences in tree topologies.

293       In many cases, the same tree topology was returned even though the number

294   of identified SNPs differed dramatically (Table 2). This result could be due to

295   multiple factors, including if and how duplicates are filtered from the reference

296   genome or other genome assemblies. With regards to NASP, erroneous SNPs

297   called in genome assemblies are likely artifacts from the whole genome

298   alignments using NUCmer. The default value for aligning through poorly scoring

299   regions before breaking an alignment in NUCmer is 200, potentially introducing

300   many spurious SNPs into the alignment, especially in misassembled regions in

301   draft genome assemblies.  By changing this value to 20, the same tree topology

302   was obtained, although many fewer SNPs (n=~100,000) were identified (Table

303   2). This value is easily altered in NASP and should be tuned based on the

304   inherent expected diversity in the chosen dataset. Additional investigation is

305   required to verify that SNPs in divergent regions are not being lost by changing

306   this parameter. Another option is to use simulated reads from the genome

307   assemblies in the SNP identification process.

308

**Phylogeny differences on the same dataset.** Previously, it has been demonstrated that different phylogenies can be obtained on the same dataset using either RAxML or FastTree2 [15]. To test this result across multiple phylogenetic inference methods, the NASP *E. coli* read dataset was used. Phylogenies were then inferred using a maximum likelihood method in RAxML, a maximum parsimony method implemented in Parsimonator, a minimum evolution method in FastTree2, and a Bayesian method implemented in Exabayes [35]. The results demonstrate variability in the placement of one genome (UMN026) depending on the method. FastTree2 and Exabayes agreed on their topologies, including 100% congruence of the replicate trees. The maximum likelihood and maximum parsimony phylogenies were slightly different (Supplemental Figure 3) and included low bootstrap replicate values at the variable node. The correct placement of UMN026 is unknown and is likely confounded by the extensive recombination observed in *E. coli* [39].

323

**Pipeline comparisons on a well characterized dataset.** To test the functionality of different SNP calling pipelines, a set of 15 finished *Yersinia pestis* genomes were compared. This set of genomes was selected because 26 SNPs in the dataset have been verified by wet-bench methods (Supplemental Table 4). Additionally, 13 known errors in the reference genome, *Y. pestis* CO92 [40], have been identified (Supplemental Table 4) and should consistently be identified in SNP discovery methods. The small number of SNPs in the dataset requires accurate SNP identification to resolve the phylogenetic relationships of these genomes.

The results demonstrate differences in the total number of SNPs called between different methods (Table 3). Most of the methods identified all 13 known sequencing errors in CO92, although Parsnp, REALPHY, and kSNPv3 failed to do so. The number of verified SNPs also varied between methods, from 21 in kSNPv3 to all 26 in multiple methods (Table 3). An analysis of wet-bench validated SNPs (n=9) that are identified in more than one genome demonstrated

339   that some methods failed to identify all of these SNPs, which could lead to a very

340   different phylogeny than the phylogeny using these SNPs that are vital for

341   resolving important phylogenetic relationships. These SNPs could represent

342   differences that could differentiate between strains in an outbreak event.

343

344   **Pipeline comparisons on a simulated set of assemblies and reads**.

345   Simulated data for *Y. pestis* were used to compare SNP identification between

346   pipelines. In this method, 3501 mutations (Supplemental Data File 2) were

347   inserted into genomes based on a published phylogeny [41] and FASTA file. Raw

348   reads were also simulated from these artificially mutated assemblies with ART to

349   generate paired end sequences. Reads and assemblies were run across all

350   pipelines, where applicable.

351   The results demonstrate that NASP identified all of the inserted SNPs using

352   raw reads, although 68 SNPs failed the proportion filter (0.90) and 232 SNPs fell

353   in duplicated regions (Table 4); some of the duplicated SNPs would also fail the

354   proportion filter. Of all other methods, only ISG identified all inserted mutations.

355   SPANDx only identified 2248 SNPs when run with default values. Parsnp

356   identified the majority of the mutations, although duplicate regions appear to

357   have also been aligned.

358   To understand how the SNPs called would affect the overall tree topology, a

359   phylogeny was inferred for each set of SNPs with RAxML. A similarity matrix was

360   made for each method based on the topological score compared to the ML

361   phylogeny inferred from the known mutations. The UPGMA dendrogram

362   demonstrates that the NASP results return a phylogeny that is more

363   representative of the "true" phylogeny than other methods (Figure 3). Without

364   removing SNPs found in duplicated regions, the NASP phylogeny was identical

365   to the phylogeny inferred from the known SNPs.

366

367

368

369

370 **Discussion**

371     Understanding relationships between bacterial isolates in a population is

372 important for applications such as source tracking, outbreak investigations,

373 phylogeography, population dynamics, and diagnostic development. With the

374 large numbers of genomes that are typically associated with these investigations,

375 methods are required to quickly and accurately identify SNPs in a reference

376 population. However, no studies have conducted a broad analysis to compare

377 published methods on real and simulated datasets to identify relevant strengths

378 and weaknesses.

379     Multiple publications have used a reference-dependent approach to identify

380 SNPs to understand population dynamics [38]. While the specific methods are

381 often published, the pipelines to run these processes are often un-published [42,

382 43], which complicates the ability to replicate results. NASP has already been

383 used to identify SNPs from multiple organisms, including fungal [44] and bacterial

384 [45, 46] pathogens. The version-controlled source code is available for NASP,

385 which should ensure the replication of results across research groups.

386     Recently it has been suggested that the use of a single reference can bias the

387 identification of SNPs, especially in divergent references [11]. In our *E. coli* test

388 set, ~29,000 fewer SNPs were called by aligning *E. coli* reads against the

389 reference genome of the outgroup, *E. fergusonii*, compared to the K-12

390 reference, although the tree topologies were identical (Table 2). In the *E. coli* test

391 set phylogeny, the major clades are delineated by enough SNPs that the loss of

392 a small percentage is insufficient to change the overall tree topology, although

393 the branch lengths were variable. In other datasets, the choice of the reference

394 should be made carefully to include as many SNPs as needed to define the

395 population structure of a given dataset.

396     According to the authors of kSNP, a k-mer-based reference-independent

397 approach, there are times where alignments are not appropriate in understanding

398 bacterial population structure [8]. In our *E. coli* analysis, reference-dependent

399 and reference-independent methods generally returned the same tree topology

400 (Table 2), with the exception of kSNPv3 and lyve-SET, using only core genome

13

401 SNPs. Using all of the SNPs identified by kSNPv3 also gave a different tree
402 topology than the other methods (Table 2). A detailed look at branch specific
403 SNPs demonstrated that using kSNP with core SNPs failed to identify most of the
404 branch specific SNPs for one of the major defining clades (Table 2). For datasets
405 that are only defined by a small number of SNPs, a method should be chosen
406 that includes as many SNPs as possible in order to maximize the relevant search
407 space. While NASP cannot truly use the pan-genome if a single reference
408 genome is chosen, it can incorporate data from all positions in the reference
409 genome if missing data are included in the alignment. A true pan-genome
410 reference can be used with NASP to more comprehensively identify SNPs, but
411 curation of the pan-genome is necessary to remove genomic elements
412 introduced by horizontal gene transfer that could potentially confound the
413 phylogeny.

414 Phylogenetics on an alignment of concatenated SNPs is thought to be less
415 preferable than an alignment that also contains monomorphic positions [11, 38].
416 However, the inclusion of monomorphic positions can drastically increase the run
417 time needed to infer a phylogeny, especially where the population structure of a
418 species can be determined by a small number of polymorphisms. Substitution
419 models are available in RAxML v8 that contain acquisition bias corrections that
420 should be considered when inferring phylogenies from concatenated SNP
421 alignments. In our *E. coli* test case, using concatenated SNPs did not change the
422 tree topology compared to a phylogeny inferred from all sites, but did affect
423 branch lengths (Supplemental Figure 2). For downstream methods that depend
424 on accurate branch lengths, decisions must be made on whether or not to
425 include monomorphic positions into the alignment. NASP provides the user with
426 the flexibility to make those decisions in a reproducible manner.

427 NASP represents a version-controlled, unit-tested pipeline for identifying
428 SNPs from datasets with diverse input types. NASP is a high throughput method
429 that can take a range of input formats, can accommodate multiple job
430 management systems, can use multiple read aligners and SNP callers, can

14

431   identify both monomorphic and polymorphic sites, and can generate core

432   genome statistics across a population.

433

434   **Figure Legends**

435   **Figure 1**. A workflow of the NASP algorithm. Optional steps are shown by

436   dashed lines.

437

438   **Figure 2**. NASP benchmark comparisons of walltime (A) and RAM (B) on a set of

439   *Escherichia coli* genomes. For the walltime comparisons, 3520 *E. coli* genomes

440   were randomly sampled ten times at different depths and run on a server with

441   856 cores. Only the matrix building step is shown, but demonstrates a linear

442   scaling with the processing of additional genomes.

443

444   **Figure 3**. Dendrogram of tree building methods on a simulated set of mutations

445   in the genome of *Yersinia pestis* Colorado 92. The topological score was

446   generated by compare2trees compared to a maximum likelihood phylogeny

447   inferred from a set of 3501 SNPs inserted by Tree2Reads. The dendrogram was

448   generated with the Neighbor method in the Phylip software package [36].

449

450   **References**

451   1.    Rasko DA, Webster DR, Sahl JW, Bashir A, Boisen N, Scheutz F, Paxinos EE,
452         Sebra R, Chen-Shan C, Iliopoulos D, et al: **Origins of the E.coli strain causing**
453         **an oubreak of hemolytic-uremic syndrome in Germany.** *N Engl J Med*
454         2011.
455   2.    Eppinger M, Mammel MK, Leclerc JE, Ravel J, Cebula TA: **Genomic anatomy**
456         **of Escherichia coli O157:H7 outbreaks.** *Proc Natl Acad Sci U S A* 2011,
457         **108:**20142-20147.
458   3.    Sahl JW, Morris CR, Emberger J, Fraser CM, Ochieng JB, Juma J, Fields B,
459         Breiman RF, Gilmour M, Nataro JP, Rasko DA: **Defining the phylogenomics**
460         **of Shigella species: a pathway to diagnostics.** *J Clin Microbiol* 2015,
461         **53:**951-960.
462   4.    Keim PS, Wagner DM: **Humans and evolutionary and ecological forces**
463         **shaped the phylogeography of recently emerged diseases.** *Nat Rev*
464         *Microbiol* 2009, **7:**813-821.
465   5.    Foster JT, Beckstrom-Sternberg SM, Pearson T, Beckstrom-Sternberg JS,
466         Chain PS, Roberto FF, Hnath J, Brettin T, Keim P: **Whole-genome-based**

| | | |
|---|---|---|
| 467 | | **phylogeny and divergence of the genus Brucella.** *J Bacteriol* 2009, |
| 468 | | **191:**2864-2870. |
| 469 | 6. | Olson ND, Lund SP, Colman RE, Foster JT, Sahl JW, Schupp JM, Keim P, |
| 470 | | Morrow JB, Salit ML, Zook JM: **Best practices for evaluating single** |
| 471 | | **nucleotide variant calling methods for microbial genomics.** *Front Genet* |
| 472 | | 2015, **6:**235. |
| 473 | 7. | Sahl JW, Beckstrom-Sternberg SM, Babic-Sternberg J, Gillece JD, Hepp CM, |
| 474 | | Auerbach RK, Tembe W, Wagner DM, Keim PS, Pearson T: **The In Silico** |
| 475 | | **Genotyper (ISG): an open-source pipeline to rapidly identify and** |
| 476 | | **annotate nucleotide variants for comparative genomcis applications.** |
| 477 | | *bioRxiv* 2015, **015578**. |
| 478 | 8. | Gardner SN, Hall BG: **When whole-genome alignments just won't work:** |
| 479 | | **kSNP v2 software for alignment-free SNP discovery and phylogenetics** |
| 480 | | **of hundreds of microbial genomes.** *PLoS ONE* 2013, |
| 481 | | **10.1371/journal.pone.0081760**. |
| 482 | 9. | Gardner SN, Slezak T: **Scalable SNP analyses of 100 bacterial or viral** |
| 483 | | **genomes.** *J Forensic Res* 2010, **1:**doi:10.4172/2157-7145.1000107. |
| 484 | 10. | Gardner SN, Slezak T, Hall BG: **kSNP3.0: SNP detection and phylogenetic** |
| 485 | | **analysis of genomes without genome alignment or reference genome.** |
| 486 | | *Bioinformatics* 2015, **31:**2877-2878. |
| 487 | 11. | Bertels F, Silander OK, Pachkov M, Rainey PB, van Nimwegen E: **Automated** |
| 488 | | **reconstruction of whole-genome phylogenies from short-sequence** |
| 489 | | **reads.** *Mol Biol Evol* 2014, **31:**1077-1088. |
| 490 | 12. | Treangen TJ, Ondov BD, Koren S, Phillippy AM: **The Harvest suite for rapid** |
| 491 | | **core-genome alignment and visualization of thousands of intraspecific** |
| 492 | | **microbial genomes.** *Genome Biol* 2014, **15:**524. |
| 493 | 13. | Sarovich DS, Price EP: **SPANDx: a genomics pipeline for comparative** |
| 494 | | **analysis of large haploid whole genome re-sequencing datasets.** *BMC Res* |
| 495 | | *Notes* 2014, **7:**618. |
| 496 | 14. | **Future directions for research on enterotoxigenic Escherichia coli** |
| 497 | | **vaccines for developing countries.** *Wkly Epidemiol Rec* 2006, **81:**97-104. |
| 498 | 15. | Pettengill JB, Luo Y, Davis S, Chen Y, Gonzalez-Escalona N, Ottesen A, Rand H, |
| 499 | | Allard MW, Strain E: **An evaluation of alternative methods for** |
| 500 | | **constructing phylogenies from whole genome sequence data: a case** |
| 501 | | **study with Salmonella.** *PeerJ* 2014, **2:**e620. |
| 502 | 16. | Delcher AL, Salzberg SL, Phillippy AM: **Using MUMmer to identify similar** |
| 503 | | **regions in large sequence sets.** *Curr Protoc Bioinformatics* 2003, **Chapter** |
| 504 | | **10:**Unit 10 13. |
| 505 | 17. | Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for** |
| 506 | | **Illumina sequence data.** *Bioinformatics* 2014, **30:**2114-2120. |
| 507 | 18. | Li H: **Aligning sequence reads, clone sequences and assembly contigs** |
| 508 | | **with BWA-MEM.** *arXivorg* 2013. |
| 509 | 19. | Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat* |
| 510 | | *Methods* 2012, **9:**357-359. |

511 20. Zaharia M, Bolosky WJ, Curtis K, Fox A, Patterson D, Shenker S, Stoica I, Karp
512     RM, Sittler T: **Faster and More Accurate Sequence Alignment with SNAP.**
513     *arXivorg* 2011, **arXiv.1111.5572 [cs.DS]**.
514 21. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis
515     G, Durbin R, Genome Project Data Processing S: **The Sequence
516     Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25:**2078-2079.
517 22. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis
518     AA, del Angel G, Rivas MA, Hanna M, et al: **A framework for variation
519     discovery and genotyping using next-generation DNA sequencing data.**
520     *Nature genetics* 2011, **43:**491-498.
521 23. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A,
522     Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA: **The Genome
523     Analysis Toolkit: a MapReduce framework for analyzing next-
524     generation DNA sequencing data.** *Genome research* 2010, **20:**1297-1303.
525 24. Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, Miller CA,
526     Mardis ER, Ding L, Wilson RK: **VarScan 2: somatic mutation and copy
527     number alteration discovery in cancer by exome sequencing.** *Genome
528     Res* 2012, **22:**568-576.
529 25. Renteria ME, Cortes A, Medland SE: **Using PLINK for Genome-Wide
530     Association Studies (GWAS) and data analysis.** *Methods Mol Biol* 2013,
531     **1019:**193-213.
532 26. Touchon M, Hoede C, Tenaillon O, Barbe V, Baeriswyl S, Bidet P, Bingen E,
533     Bonacorsi S, Bouchier C, Bouvet O, et al: **Organised genome dynamics in
534     the Escherichia coli species results in highly diverse adaptive paths.**
535     *PLoS Genet* 2009, **5:**e1000344.
536 27. Huang W, Li L, Myers JR, Marth GT: **ART: a next-generation sequencing
537     read simulator.** *Bioinformatics* 2012, **28:**593-594.
538 28. Blattner FR, Plunkett G, 3rd, Bloch CA, Perna NT, Burland V, Riley M, Collado-
539     Vides J, Glasner JD, Rode CK, Mayhew GF, et al: **The complete genome
540     sequence of Escherichia coli K-12.** *Science* 1997, **277:**1453-1462.
541 29. Cui Y, Yu C, Yan Y, Li D, Li Y, Jombart T, Weinert LA, Wang Z, Guo Z, Xu L, et al:
542     **Historical variations in mutation rate in an epidemic pathogen, Yersinia
543     pestis.** *Proc Natl Acad Sci U S A* 2013, **110:**577-582.
544 30. Angiuoli SV, Salzberg SL: **Mugsy: Fast multiple alignment of closely
545     related whole genomes.** *Bioinformatics* 2010.
546 31. Katz LS, Petkau A, Beaulaurier J, Tyler S, Antonova ES, Turnsek MA, Guo Y,
547     Wang S, Paxinos EE, Orata F, et al: **Evolutionary dynamics of Vibrio
548     cholerae O1 following a single-source introduction to Haiti.** *MBio* 2013,
549     **4**.
550 32. Sahl JW, Steinsland H, Redman JC, Angiuoli SV, Nataro JP, Sommerfelt H,
551     Rasko DA: **A comparative genomic analysis of diverse clonal types of
552     enterotoxigenic Escherichia coli reveals pathovar-specific conservation.**
553     *Infect Immun* 2011, **79:**950-960.
554 33. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and
555     post-analysis of large phylogenies.** *Bioinformatics* 2014.

556  34.  Price MN, Dehal PS, Arkin AP: **FastTree 2--approximately maximum-likelihood trees for large alignments.** *PLoS One* 2010, **5:**e9490.
558  35.  Aberer AJ, Kobert K, Stamatakis A: **ExaBayes: massively parallel bayesian tree inference for the whole-genome era.** *Mol Biol Evol* 2014, **31:**2553-2556.
561  36.  Felsenstein J: **PHYLIP (Phylogeny Inference Package) version 3.6.** 3.6 edition. University of Washington, Seattle: Department of Genome Sciences; 2005.
564  37.  Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM: **A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3.** *Fly (Austin)* 2012, **6:**80-92.
568  38.  Leache AD, Banbury BL, Felsenstein J, de Oca AN, Stamatakis A: **Short Tree, Long Tree, Right Tree, Wrong Tree: New Acquisition Bias Corrections for Inferring SNP Phylogenies.** *Syst Biol* 2015.
571  39.  Dykhuizen DE, Green L: **Recombination in Escherichia coli and the definition of biological species.** *Journal of bacteriology* 1991, **173:**7257-7268.
574  40.  Parkhill J, Wren BW, Thomson NR, Titball RW, Holden MT, Prentice MB, Sebaihia M, James KD, Churcher C, Mungall KL, et al: **Genome sequence of Yersinia pestis, the causative agent of plague.** *Nature* 2001, **413:**523-527.
577  41.  Nawrocki EP, Kolbe DL, Eddy SR: **Infernal 1.0: inference of RNA alignments.** *Bioinformatics* 2009, **25:**1335-1337.
579  42.  den Bakker HC, Allard MW, Bopp D, Brown EW, Fontana J, Iqbal Z, Kinney A, Limberger R, Musser KA, Shudt M, et al: **Rapid whole-genome sequencing for surveillance of Salmonella enterica serovar enteritidis.** *Emerg Infect Dis* 2014, **20:**1306-1314.
583  43.  Hsu LY, Harris SR, Chlebowicz MA, Lindsay JA, Koh TH, Krishnan P, Tan TY, Hon PY, Grubb WB, Bentley SD, et al: **Evolutionary dynamics of methicillin-resistant Staphylococcus aureus within a healthcare system.** *Genome Biol* 2015, **16:**81.
587  44.  Engelthaler DM, Hicks ND, Gillece JD, Roe CC, Schupp JM, Driebe EM, Gilgado F, Carriconde F, Trilles L, Firacative C, et al: **Cryptococcus gattii in North American Pacific Northwest: whole-population genome analysis provides insights into species evolution and dispersal.** *MBio* 2014, **5:**e01464-01414.
592  45.  Sahl JW, Schupp JM, Rasko DA, Colman RE, Foster JT, Keim P: **Phylogenetically typing bacterial strains from partial SNP genotypes observed from direct sequencing of clinical specimen metagenomic data.** *Genome Med* 2015, **7:**52.
596  46.  Sahl JW, Sistrunk JR, Fraser CM, Hine E, Baby N, Begum Y, Luo Q, Sheikh A, Qadri F, Fleckenstein JM, Rasko DA: **Examination of the Enterotoxigenic Escherichia coli Population Structure during Human Infection.** *MBio* 2015, **6:**e00501.
600

**Table 1**. An overview of commonly used SNP pipelines

| Pipeline name | Supported data types | Output type | Parallel Job management support? |
|---|---|---|---|
| NASP | FASTA,BAM,SAM,VCF,FASTQ,FASTQ.GZ | matrix, VCF, FASTA | SGE,SLURM,TORQUE |
| ISG | FASTA,BAM,VCF,FASTQ,FASTQ.GZ | matrix, FASTA | No |
| Parsnp | FASTA | gingr file, phylogeny, FASTA, VCF | No |
| REALPHY | FASTA*, FASTQ, FASTQ.GZ | multi-FASTA, phylogeny | No |
| SPANDx | FASTQ.GZ | Nexus file | SGE,SLURM,TORQUE |
| CFSAN | FASTQ, FASTQ.GZ | SNP list, FASTA | SGE,TORQUE |
| kSNPv3 | FASTA | Matrix, FASTA, phylogeny | No |
| Mugsy | FASTA | MAF file | No |
| LYVE-set | FASTQ.GZ, FASTA* | matrix, FASTA, phylogeny | SGE |

601    *generates simulated reads

**Table 2. SNP calling results on a set of 21 *E. coli* genomes**

| Method | Reference | Data Type | Parameters | #SNPs considered | #Total sites | Walltime (single node - 8 cores) | Topological score | #defining SNPs |
|---|---|---|---|---|---|---|---|---|
| NASP | K12 MG1655 | Assemblies | Default | 267978[a] | 2322434 | 10m00s | 100% | 809 |
| NASP | K12 MG1655 | Assemblies | NUCmer (-b 20) | 162758[a] | 1839583 | 10m00s | 100% | 744 |
| NASP | K12 MG1655 | ART PE reads | BWA, GATK,MinDepth=3,MinAF=0.90 | 170208[a] | 1984510 | 1h43m00s | 100% | 826 |
| NASP | *E. fergusonii* 35469 | Assemblies | Default | 244262[a] | 2227038 | 10m00s | 100% | 741 |
| NASP | *E. fergusonii* 35469 | ART PE reads | BWA,GATK,MinDepth=3,MinAF=0.90 | 141238[a] | 1813349 | 1h17m10s | 100% | 748 |
| ISG | K12 MG1655 | Assemblies | Default | 268524[a] | N/A | 6m47s | 100% | 810 |
| ISG | K12 MG1655 | ART PE reads | minaf 0.9, mindp 3 | 206193[a] | N/A | 14m45s | 100% | 824 |
| Parsnp | K12 MG1655 | Assemblies | "-c d" | 151256[a] | 1682404 | 4m35s | 100% | 777 |
| REALPHY | K12 MG1655 | REALPHY SE reads | Default | 171828[a] | 1897146 | 3h11m00s | 100% | 779 |
| kSNPv3 | N/A | Assemblies | -core | 20587[a] | N/A | 27m58s | 91.80% | 5 |
| kSNPv3 | N/A | Assemblies | Default | 284134 | N/A | 27m58s | 95.80% | 547 |
| Spandx | K12 MG1655 | ART PE reads | Default | 95214[a] | N/A | N/A | 100% | 709 |
| CFSAN | K12 MG1655 | ART PE reads | Default | 128512[a] | N/A | 1h56m00s | 100% | 808 |
| Mugsy | N/A | Assemblies | Default | 307072[a] | 2478794 | 1h39m03s | 100% | unknown |
| lyve-SET | K12 MG1655 | ART PE reads | min_coverate 3, min_alt_frac 0.9 | 163118[a] | 1183153 | 6h25m | 85% | 329 |

602 [a] strictly core genome SNPs

603

604

605

606

607

608

609
610
611
612
613

**Table 3. SNP calling results on a set of *Y. pestis* genomes**

| Method | Data type | Parameters | #called SNPs | #CO92 errors (n=13) | #verified SNPs (n=26) | Vital SNP (n=9) |
|---|---|---|---|---|---|---|
| NASP | ART simulated reads | BWA,GATK,MinDepth=3,MinAF=0.90 | 147 | 13 | 26 | 9 |
| NASP | assemblies | default | 181 | 13 | 26 | 9 |
| ISG | ART simulated reads | minaf=3, mindp = 0.9 | 151 | 13 | 26 | 9 |
| ISG | assemblies | default | 177 | 13 | 26 | 9 |
| Parsnp | assemblies | default | 141 | 12 | 23 | 7 |
| REALPHY | REALPHY simulated reads | default | 163 | 12 | 25 | 9 |
| SPANDx | ART simulated reads | default | 150 | 13 | 25 | 9 |
| kSNPv3 | assemblies | k=19 | 130 | 11 | 21 | 5 |
| CFSAN | ART simulated reads | default | 250 | 13 | 26 | 9 |
| lyve-SET | ART simulated reads | min_coverage 3, min_alt_frac 0.9 | 402 | 13 | 26 | 9 |

614
615
616

20

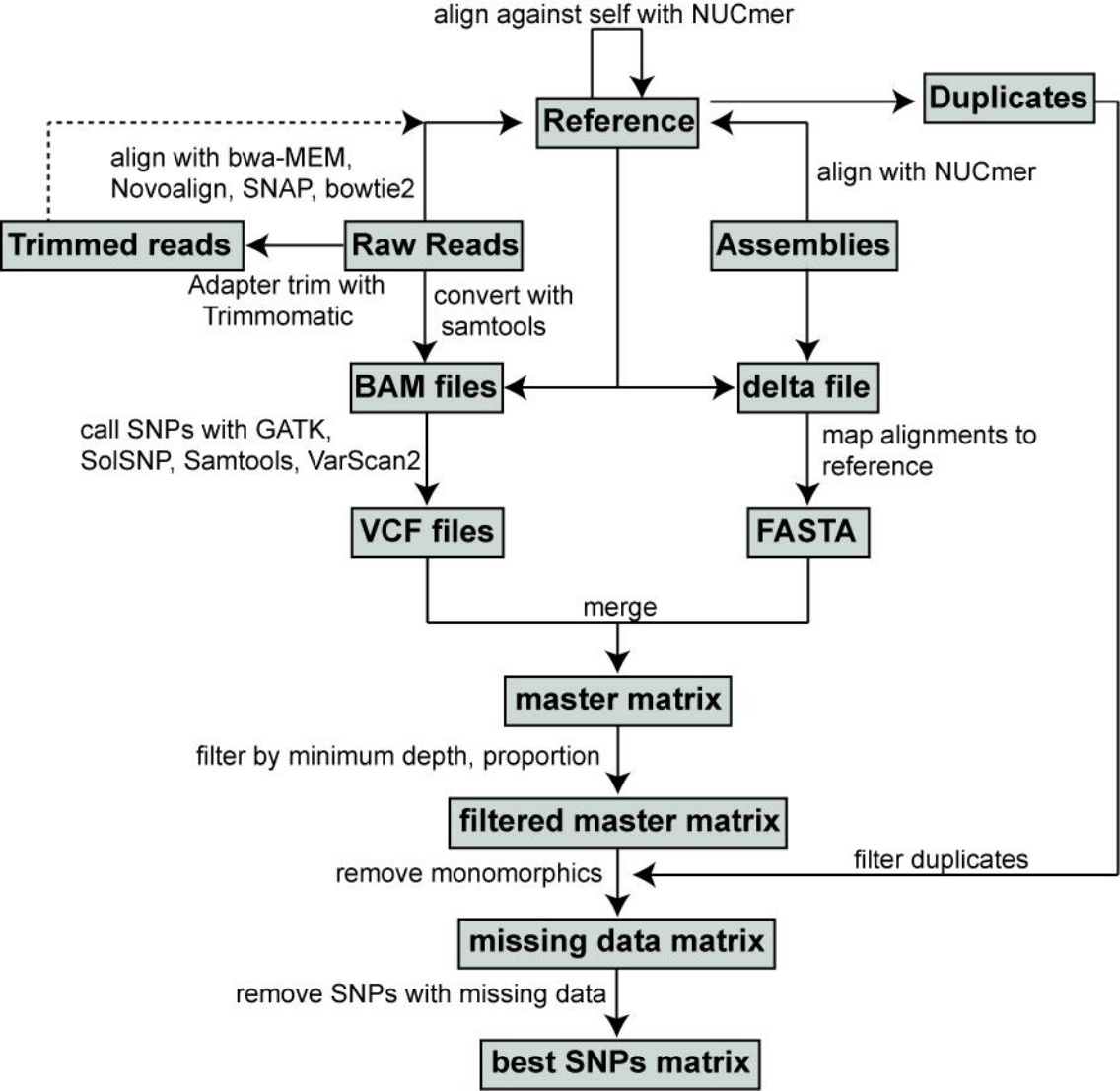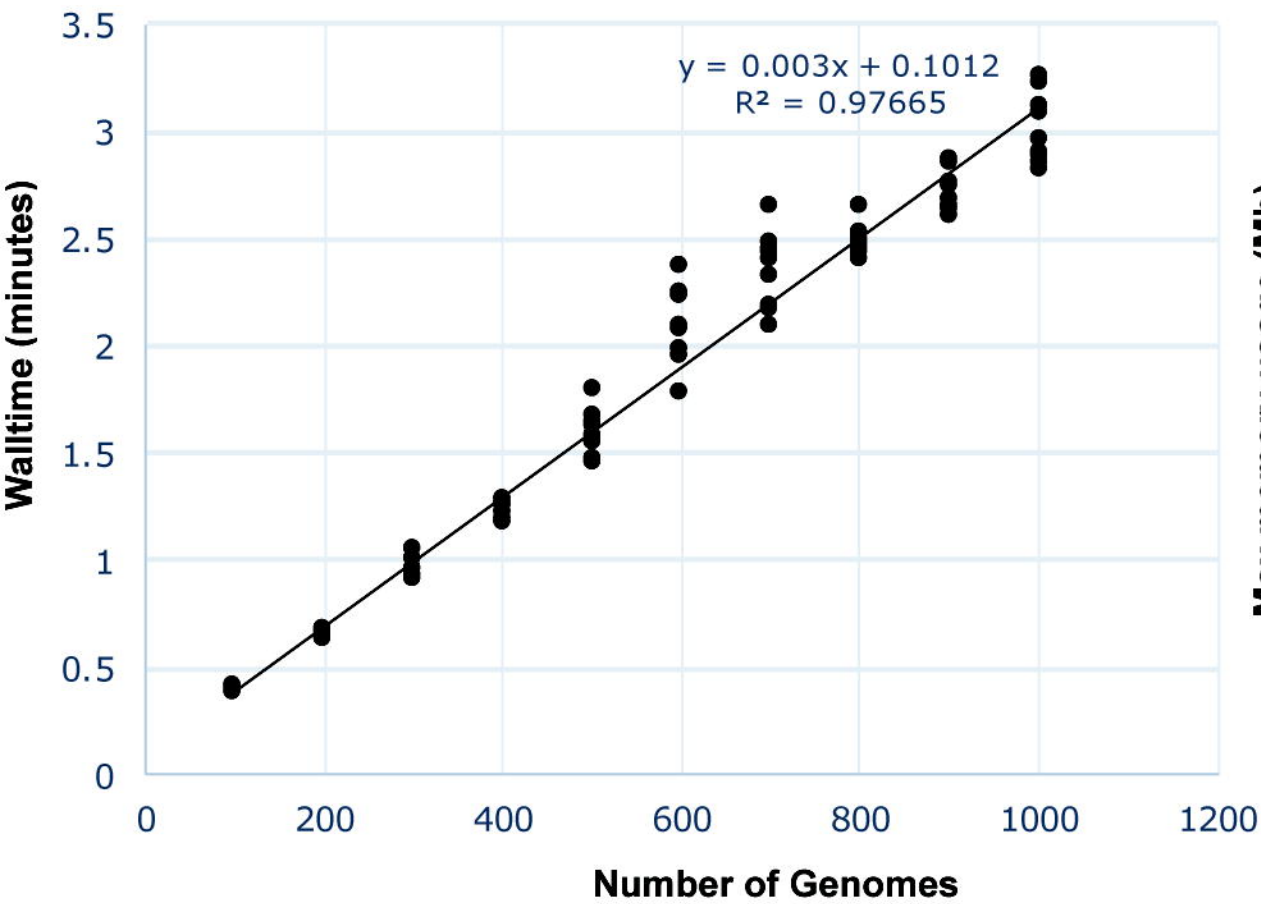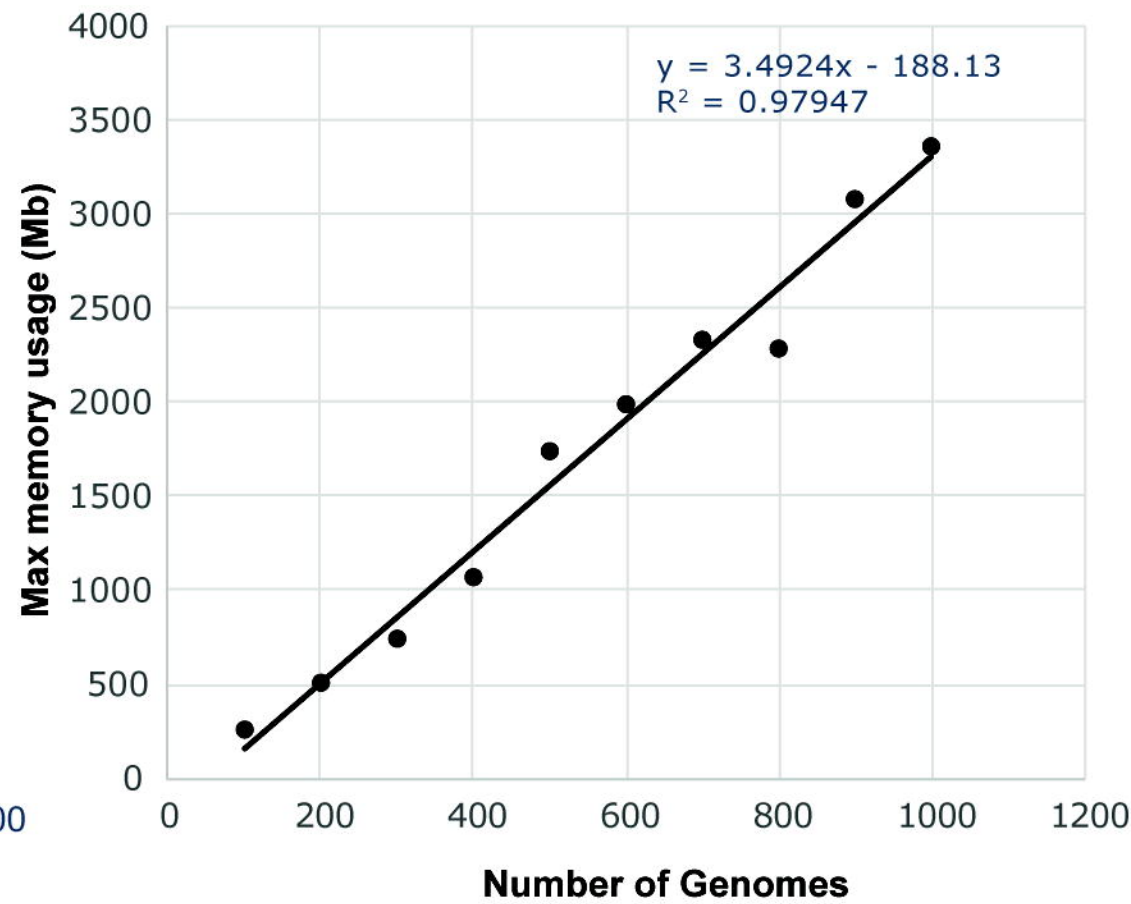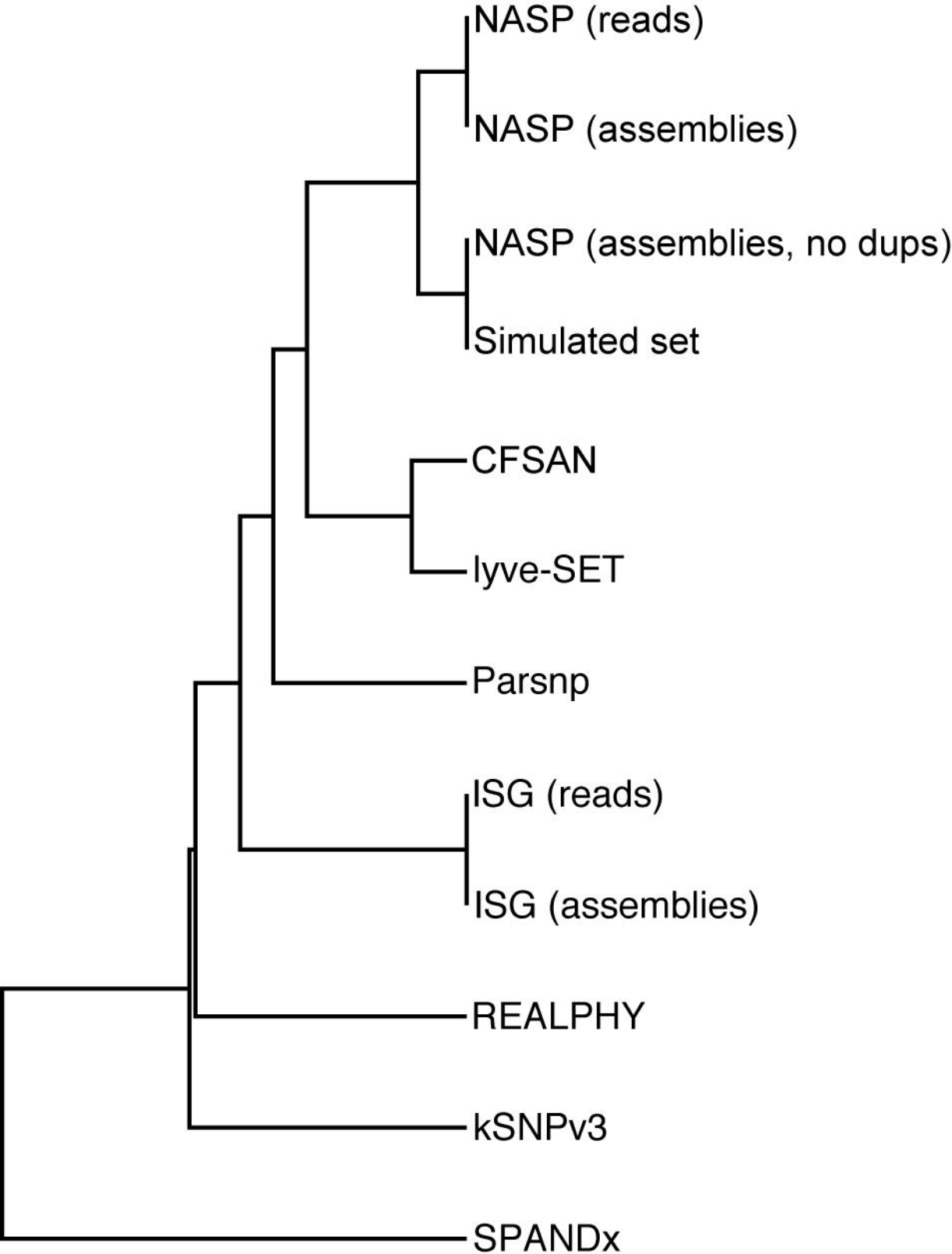**Table 4**. Simulated data results

| Method | Data type | #called SNPs | SNPs in duplicated regions | Filtered SNPs | Total SNPs | Topologica score |
|--------|-----------|--------------|----------------------------|---------------|------------|-------------------|
| NASP | simulated reads | 3202 | 232 | 67 | 3501 | 98.50% |
| NASP | simulated assemblies | 3269 | 232 | N/A | 3501 | 98.50% |
| Parsnp | simulated assemblies | 3492 | unknown | N/A | 3492 | 95.60% |
| ISG | simulated reads | 3258 | 126 | 8 | 3392 | 92.40% |
| ISG | simulated assemblies | 3266 | 235 | N/A | 3501 | 95.60% |
| SPANDx | simulated reads | 2132 | unknown | 116 | 2248 | 87.20% |
| CFSAN | simulated reads | 3290 | unknown | unknown | 3290 | 95.30% |
| REALPHY | simulated assemblies | 3320 | unknown | unknown | 3320 | 91.60% |
| kSNPv3 | simulated assemblies | 3304 | unknown | N/A | 3304 | 91.90% |
| lyve-SET | simulated reads | 3460 | unknown | unknown | 3460 | 95.80% |

617
618
619
620
621
622
623
624
625
626
627
628
629

21

align against self with NUCmer

**Reference**

**Duplicates**

align with bwa-MEM, Novoalign, SNAP, bowtie2

align with NUCmer

**Trimmed reads**

**Raw Reads**

**Assemblies**

Adapter trim with Trimmomatic

convert with samtools

**BAM files**

**delta file**

call SNPs with GATK, SolSNP, Samtools, VarScan2

map alignments to reference

**VCF files**

**FASTA**

merge

**master matrix**

filter by minimum depth, proportion

**filtered master matrix**

remove monomorphics

filter duplicates

**missing data matrix**

remove SNPs with missing data

**best SNPs matrix**

**A. Elapsed time**

$y = 0.003x + 0.1012$
$R^2 = 0.97665$

Walltime (minutes)

Number of Genomes

**B. RAM usage**

$y = 3.4924x - 188.13$
$R^2 = 0.97947$

Max memory usage (Mb)

Number of Genomes

## All positions

32449 ─── *E. fergusonii* 35469

8055
15748 ─── IAI39
16386
5843.5 ─ 536
2225.5
6751.5 ─ CFT073
6883.5 ─ ED1A
4403.5
139 ─ UTI89
4766.5
1957 ─ S88
190
258 ─ APEC_01

32449

9223
14579 ─── UMN026
5443.5
13947 ─── *S. dysenteriae* 197
10914
26 ─ O157:H7_EDL933
26 ─ O157:H7_sakai

9379.5
12527.5 ── K12_W3110
K12_MG1655
Reference
5482
842.5 ─ *S. flexneri* 5B
10034
184 ─ *S. flexneri* 2A
781.5
282 ─ *S. flexneri* 2A 301
4042
3554
5360.5 ─ IAI1
5527.5 ─ 55989
4980
8882 ─── *S. sonnei* 046
2081
9164 ─── *S. boydii* 227

├────── 20000 SNPs ──────┤

## Polymorphisms only

32342.5 ─── *E. fergusonii* 35469

8027
15699 ─── IAI39
16338
5809 ─ 536
2223.5
6723.5 ─ CFT073
6858.5 ─ ED1A
4384
139 ─ UTI89
4763.5
1947 ─ S88
190
257 ─ APEC_01

32342.5

9180
14527 ─── UMN026
5432.5
13896 ─── *S. dysenteriae* 197
10876
26 ─ O157:H7_EDL933
26 ─ O157:H7_sakai

9349.5
12481.5 ── K12_W3110
K12_MG1655
Reference
5475
838.5 ─ *S. flexneri* 5B
9984
178 ─ *S. flexneri* 2A
778.5
275 ─ *S. flexneri* 2A 301
4029
3541
5353.5 ─ IAI1
5516.5 ─ 55989
4959
8845 ─── *S. sonnei* 046
2065
9137 ─── *S. boydii* 227

├────── 20000 SNPs ──────┤

● Boostrap or Posterior probabilities = 100