1    **From raw reads to trees: Whole genome SNP phylogenetics across the tree of**

2    **life**

3

4    Sanaa A. Ahmed[1], Chien-Chi Lo[1], Po-E Li[1], Karen W. Davenport[1], and Patrick S. G. Chain[1]

5

6    [1]Biome Sciences, Bioscience Division, Los Alamos National Laboratory, MS-M888, Los

7    Alamos, NM 87545

8    **Corresponding author:**

9    Patrick S.G. Chain
10   Scientist IV
11   Los Alamos National Laboratory
12   MS-M888
13   Los Alamos, NM 87545
14   Email: pchain@lanl.gov
15   Tel: 505-665-4019
16   Cell: 505-500-6072

17   **Running Title:  Raw sequence reads to whole genome SNP phylogenies**

18   **Keywords:** Whole-genome SNP phylogenies; phylogenomics; positive selection; Illumina raw
19   data analysis; metagenomic reads to phylogenetic trees; PhaME

20

21

22

23

24

25

26

27

28

29

30 **ABSTRACT**

31 Next-generation sequencing is increasingly being used to examine closely related organisms.

32 However, while genome-wide single nucleotide polymorphisms (SNPs) provide an excellent

33 resource for phylogenetic reconstruction, to date evolutionary analyses have been performed

34 using different ad hoc methods that are not often widely applicable across different projects. To

35 facilitate the construction of robust phylogenies, we have developed a method for genome-wide

36 identification/characterization of SNPs from sequencing reads and genome assemblies. Our

37 phylogenetic and molecular evolutionary (PhaME) analysis software is unique in its ability to

38 take reads and draft/complete genome(s) as input, derive core genome alignments, identify

39 SNPs, construct phylogenies and perform evolutionary analyses. Several examples using

40 genomes and read datasets for bacterial, eukaryotic and viral linages demonstrate the broad and

41 robust functionality of PhaME. Furthermore, the ability to incorporate raw metagenomic reads

42 from clinical samples with suspected infectious agents shows promise for the rapid phylogenetic

43 characterization of pathogens within complex samples.

44

45

46

47

48

49

50

51

52

53

54

55

56 **INTRODUCTION**

57 The reconstruction of evolutionary history using phylogenetics has been a fundamental method

58 applied to many areas of biology. Single nucleotide polymorphisms (SNPs), one of the dominant

59 forms of evolutionary change, have become an indispensable tool for phylogenetic analyses.

60 Phylogenies in the pre-genomic era relied on SNPs found within single genes as the evolutionary

61 signal, and later incorporated multiple loci, such as in multiple locus sequence typing (MLST).

62 Although still valuable, these methods are limited to differences found in short sequences

63 representing only a fraction of the genome, and are unable to capture the complete variation

64 within species. As a result, these methods generally provide a low or an insufficient phylogenetic

65 signal as gene-based trees do not always reflect the true species tree(Pamilo and Nei 1988).

66 While phylogenetic analyses that use many orthologs are an improvement, multi-ortholog

67 methods only utilize variable sites within annotated coding regions that pass a user-specified

68 orthology test, and, more problematically, require both an accurate and existing annotation of all

69 of the orthologous genes found within assembled genomes.

70

71 Whole-genome SNPs are one of the most commonly used features for measuring phylogenetic

72 diversity since they can identify key phylogenetic clades for most organisms and can help

73 resolve both short and long branches(Girault et al. 2014; Griffing et al. 2015). Additionally, since

74 selectively neutral SNPs accumulate at a uniform rate, they can be used to measure divergence

75 between species, as well as strains(Schork et al. 2000; Filliol et al. 2006). Furthermore, due to

76 the large number of SNPs found along the length of entire genomes, the use of whole-genome

77 SNPs minimizes the impact of sequencing and assembly errors, as well as individual genes under

78 strong selective pressure.

79

80 Although genome-wide sequencing now allows examination of the entire genome, finished

81 genomes are only rarely produced due to platform biases, or computational and cost limitations

82 (e.g. 111,857 assembled genomes versus 7,555,153 SRA datasets in NCBI). However,

83 comparative genomics, including SNP- and ortholog gene-based phylogenetic analyses usually

84 require assembled or finished genomes(Ben-Ari et al. 2005; Foster et al. 2009). Here we present

85 PhaME, a whole-genome SNP-based approach that can make use of available completed

3

86  genomes, draft assemblies (contigs), and raw reads to perform Phylogenetic and Molecular

87  Evolutionary analyses.

88

89  Several methods for whole-genome SNP discovery or phylogenetics have been described:

90  SNPsFinder(Song et al. 2005), PhyloSNP(Faison et al. 2014), kSNP(Gardner and Hall 2013),

91  WG-FAST(Sahl et al. 2015), and CFSAN(Davis et al. 2015). SNPsFinder requires assembled

92  genomes, uses the time-consuming megaBLAST program, and only provides a table of identified

93  SNPs. PhyloSNP requires pre-aligned data (such as vcf files or tab-delimited lists of SNPs), does

94  not scale to large datasets, and only produces maximum parsimony trees based on the presence

95  or absence of SNPs. Like PhaME, kSNP, WG-FAST, and CFSAN can analyze raw reads to

96  identify a core genome (the conserved portion among all genomes). However kSNP is restricted

97  to finding central SNPs within an optimal kmer window, the size of which must first be provided

98  by the user and influences the resulting tree. WG-FAST is a rapid method that identifies the most

99  closely related known genome for given input data, but requires a pre-formatted SNP matrix of

100 known SNP positions for a target organism and a pre-computed phylogeny. CFSAN aligns read

101 datasets against a designated reference genome to generate a SNP matrix, but only allows the

102 inclusion of a single reference genome and does not build phylogenies. Because each of these

103 pipelines perform only part of the steps required to infer phylogenies from a particular form of

104 sequencing data (i.e. either reads or contigs), we have developed an integrated PhaME analysis

105 pipeline which provides rapid tree construction from assemblies and reads and downstream

106 evolutionary analysis of functional genes, with the ability to incorporate pre-constructed

107 alignments and phylogenies.

108 PhaME is an open source utility (https://github.com/LANL-Bioinformatics/PhaME) that

109 combines algorithms for whole-genome alignments, read mapping, and phylogenetic

110 construction, and uses in-house scripts to discover the core genome and core SNPs, infer trees,

111 and conduct further molecular evolutionary analyses. Because of the algorithm's design, PhaME

112 can also provide phylogenetic analyses of a target organism present even at low abundance

113 within metagenome samples, a feature unique to PhaME. This method is rapid and extensible,

114 able to build upon an existing multiple sequence alignment for any set of related genomes and

115 infer a new tree. Similarly, trees for subsets of the input genomes can readily be calculated. We

4

116    demonstrate PhaME's ability to construct robust phylogenies that include raw sequencing read

117    datasets as input (including metagenome samples) by building phylogenies with up to 560

118    genomes and spanning the tree of life.

119

120    **RESULTS**

121    **Implementation of PhaME across the tree of life**

122    PhaME was designed to identify SNPs anywhere in the genomes of closely related organisms

123    and to infer phylogenetic trees using any type of sequence input, including finished genomes,

124    assembled contigs of draft genomes, or even FASTQ raw read datasets (see Methods). We tested

125    PhaME on: 1) two large groups of bacteria, 2) the largest available dataset for a eukaryotic genus

126    (*Saccharomyces*), and 3) a viral dataset that includes many assemblies from the recent 2014

127    *Zaire ebolavirus* outbreak. We have further examined the robustness of how PhaME handles raw

128    read datasets, including those from metagenomic samples that contain varied amounts of the

129    target organism. Table 1 summarizes the number of genomes and raw reads used in the analyses,

130    the average genome size per group, the calculated core genome sizes (and percent of the average

131    genome size), the number of SNPs found in the core genome (and percent of the core genome

132    size), and the number of coding (CDS) SNPs (and percent of the total SNPs), as well as the

133    runtime information.

134

135    *Fine-scale phylogenetic analysis among closely related bacterial genera*

136    The model bacterium *Escherichia coli* has been extensively studied, including its diversity and

137    phylogenetic history. In previous studies, phylogenetic analyses using either 31 essential

138    genes(Sahl et al. 2012) or all identified core genes(Ahmed et al. 2012) have shown that 1) the

139    named *E. coli* species can be clustered into five main phylogenetic groups (A, B1, B2, D, E), 2)

140    *Shigella sonnei, Shigella boydii,* and *Shigella flexneri* form two phylogenetic groups (groups SF

141    and SS) within the *E. coli* phylogeny, and 3) *Shigella dysenteriae* clusters with the *E. coli*

142    phylogenetic group E. These analyses demonstrate that the *Shigella* species are closely related to

143    *E. coli* and not truly a separate genus(Chaudhuri and Henderson 2012), and that other

144    *Escherichia* species (E.g. *E. blattae*) are likely inaccurately classified as *Escherichia*(Priest and

145    Barker 2010). More recently, environmental organisms assigned to *Escherichia* appear to form

146    independent lineages, or 'cryptic clades' and are also closely related to *Salmonella*(Walk et al.

147    2009). While these studies have illustrated the diversity of *Escherichia*, each study has imposed

148    a different method of phylogenetic analysis sometimes resulting in incongruent trees(Walk et al.

149    2009; Luo et al. 2011), and none of the studies incorporates strains of all lineages or are able to

150    provide a high-resolution phylogeny. We used this well-characterized group to both validate

151    PhaME and to help better understand the phylogenetic relationships among members of this

152    group at a more granular level than previous studies.

153

154    Using a total of 413 NCBI complete and draft enteric genomes (including *Shigella*, *Escherichia*

155    and *Salmonella*; Supplemental Table S1), PhaME calculated a core genome of 168,191 bp with

156    49,337 SNP positions (of which 47,813 resided within annotated coding sequences; Table 1) and

157    inferred a phylogenetic tree (Figure 2; Supplemental Fig. S1). Despite the inclusion of distantly

158    related organisms, the phylogeny shows an identical topology for the *E. coli* and *Shigella* clades

159    compared with prior studies(Ahmed et al. 2012; Sahl et al. 2012). Additionally, this tree resolves

160    the evolutionary relationships among the environmental cryptic *Escherichia* lineages. For

161    example, consistent with the 2009 MLST study, but in contrast with the 2011 single copy core

162    gene study, the *E. albertii* lineage diverged before *E. fergusonii*(Walk et al. 2009; Priest and

163    Barker 2010). In addition, this tree also supports the renaming of *E. blattae* to *Shimwellia*

164    *blattae*(Walk et al. 2009), and further suggests that *E. hermanii* should also be reclassified into a

165    different genus.

166

167    While this inter-genus tree helped determine the relationships among these divergent lineages,

168    the available SNP matrix (see Methods) allows PhaME to rapidly perform a more refined

169    analysis of any genome subset, for example, just the *Escherichia* genus (Supplemental Fig. S2),

170    or just the *E. coli* clade (Supplemental Fig. S3). The core genome for all *Escherichia* (excluding

171    the renamed *S. blattae* and *E. hermanii*) was calculated to be 1,468,011 bp, with 461,639 SNP

172    positions, roughly ten times larger than the core (and core SNPs) found when including the other

173    genera. The core genome for all *E. coli* (and *Shigella*) genomes is 1,823,862 bp with 541,834

174    SNP positions, providing even more resolution within those lineages. While the improved

175    resolution of subtrees likely provides more accurate estimates of lineage evolution (branch

176    lengths), no major topological differences were observed when comparing the original tree with

177    the subtrees, supporting the robustness of PhaME even across genera.

6

178

179　*Establishing robust placement of read datasets within a phylogeny*

180　The *Burkholderia* genus is a highly diverse group, occupying a wide range of ecological niches,

181　and can be either free-living or symbionts. Some members are commensals or pathogens of a

182　variety of eukaryotic organisms. On a genome level, *Burkholderia* species have either 2 or 3

183　designated chromosomes whose sizes can vary tremendously. While the smallest genome, for *B.*

184　*rhixozinica,* has 1 chromosome and 1 megaplasmid with a total genome size of 3.6 Mbp(Lackner

185　et al. 2011), *B. xenovorans* has a much larger genome of 9.8 Mbp spanning 2 chromosomes and

186　a megaplasmid(Daligault et al. 2014). As of 2009, there are over 30 described species of

187　*Burkholderia*(Coenye and Vandamme 2003), many of which have a sequenced representative,

188　(the biothreat pathogens *B. mallei* and *B. pseudomallei* are overrepresented among the sequenced

189　genomes). One group commonly known as the *B. cepacia* complex (*Bc*c) was first described as a

190　cluster of 9 closely related species and now contains 15 distinct species. Over the years, the *Bc*c

191　have been uncharacteristically difficult to discriminate using conventional polyphasic, 16S, or

192　MLST approaches, and there are still longstanding issues in resolving their evolutionary

193　history(Coenye and Vandamme 2003; Storms et al. 2004; Reik et al. 2005; Baldwin et al. 2007;

194　Vanlaere et al. 2009).

195

196　We used PhaME to infer a genus-level phylogenetic tree (Figure 3; Supplemental Fig. S4) using

197　163 complete or draft genomes as well as 31 read files (totaling 194 datasets and 0.783 TB of

198　information; Supplemental Table S2). PhaME calculated a core genome for the *Burkholderia*

199　genus of 1,028,251 bp with 72,968 total SNPs (of which 72,642 are found in coding sequences).

200　In comparison with the *Escherichia* genus (excluding *Shimwella blattae* and *E. hermanii*), this

201　core is smaller with respect to the average genome size and the number of SNPs is proportionally

202　even smaller. This confirms that the *Burkholderia* are highly diverse with a large accessory

203　genome(Chain et al. 2006), but that the smaller core genome itself is highly conserved.

204

205　The phylogenetic tree (Figure 3; Supplemental Fig. S4) highlights the recent clonal derivation of

206　*B. mallei* from *B. pseudomallei*, with *B. pseudomallei* 576 as the most closely related genome,

207　and recapitulates the paraphyletic nature of the *B. pseudomallei* strains(Godoy et al. 2003) when

208　*B. mallei* is considered its own species.  This tree also shows two well-supported (bootstrap

209   value of 99) separate clades representing the *B. thailandensis* genomes suggesting a possible

210   reclassification of one of these two groups as a separate species. In addition, while *B. gladioli*

211   has been described as closely related to the *Bcc* species(Brisse et al. 2000; Coenye et al. 2001),

212   not only is it monophyletic with *B. glumae*, but together these two species are found as an

213   outgroup to the rest of the *Bcc* genomes and the *B. pseudomallei* group (further supported by the

214   subtree in Supplemental Fig. S5). The placement of the endosymbiont *B. rhizoxinica* as an

215   outgroup with respect to the environmental isolates was determined with PhaME when including

216   three genomes of *Rastonia pickettii* as outgroup (Supplemental Fig. S6).

217

218   The ability to select a subset of genomes for analysis can provide insight into not only the

219   consistency of the topology with the original tree, but can show differences in the resulting core

220   size and the SNPs within the core. For example, while the topology of the *Bcc* complex subtrees

221   remains the same as the original tree, the core genome size increases to 1,800,938 bp with

222   391,765 SNPs (380,062 within CDS). However, when examining only the *B. pseudomallei*

223   (including *B. mallei*) lineage, the core genome is even larger (3,351,112 bp) yet the number of

224   core SNPs is much smaller (89,399 with only 71,479 within CDS). This is in sharp contrast with

225   the *Escherichia* and the *E. coli* core genomes and SNP counts, where an increase in core genome

226   size is accompanied by an equivalent increase in SNPs. Furthermore, the number of SNPs used

227   to infer the *Burkholderia* genus tree and the *B. pseudomallei* tree are similar, yet the *B.*

228   *pseudomallei* tree provides much more discriminatory power among these closely related strains.

229   Given these patterns, we reason that the nucleotides at SNP positions used for the *Burkholderia*

230   genus tree are in fact identical among the *B. pseudomallei* genomes, and are not the same as

231   those used for *B. pseudomallei* tree.

232

233   Because one of the other unique features of PhaME, beyond rapid subtree inference, is that it

234   allows the inclusion of raw read datasets into a whole genome SNP phylogeny, we evaluated the

235   accuracy of placement of read datasets compared with the resulting assemblies and finished

236   genomes. Both the whole genus tree (Supplemental Fig. S4) and a more detailed *B.*

237   *pseudomallei*/*B. mallei* subtree (Supplemental Fig. S7) correctly place the reads with the contigs

238   and genomes that were derived from those reads (phylogenies using only the reads and not the

239   contigs/genomes resulted in consistent topologies, data not shown). This self-consistency

240    supports the ability of PhaME to conduct robust phylogenetic analysis without the need for

241    assembling raw sequencing data prior to inclusion in a tree.

242

243    ***Application of PhaME to eukaryotic genomes***

244    Because PhaME should be readily applied to any other group of closely related genomes, we

245    decided to implement tests beyond bacterial lineages. Fungi are known to be a difficult group in

246    terms of phylogenetic analysis(James et al. 2006a; James et al. 2006b; Hibbett et al. 2007). The

247    phylogenetic placement of fungal species display disparities between trees based on gene

248    sequence analyses and those based on morphological characteristics (such as modes of

249    reproduction). This is especially true of the '*Saccharomyces* complex', where 18S and 26S

250    rDNA comparisons do not show well-supported clades, resolving only the most closely related

251    species(Kurtzman and Robnett 2003).

252

253    Due to the complexity of assembling eukaryotic genomes, there are few reference draft or

254    complete genome assemblies for eukaryotic species. This is the case for the *Saccharomyces*

255    genus, which only has a single complete reference genome. For eukaryotic genomes, the ability

256    to make use of raw reads can be of great value in characterizing those genomes. A total of 79

257    genome projects consisting of 1 complete genome for *S. cerevisiae S288c* (16 chromosomes that

258    total 12,071,326 bp in length), 76 draft genomes, and 2 raw read sets (Supplemental Table S3)

259    were used as input to the PhaME pipeline. All 16 chromosomes of the yeast genome were used

260    for this analysis and, using all the input datasets, we generated the first whole genome phylogeny

261    using all available *Saccharomyces* genomes. This phylogenetic tree is based on only 463 whole-

262    genome SNP positions (explaining some of the low bootstrap confidence values in the tree), all

263    of which reside within CDS regions identified in 1,141,335 bp of core genome (Supplemental

264    Fig. S8). This also highlights the highly conserved nature of a eukaryote's core genome. A

265    refined analysis focusing solely on the *S. cerevisiae* clade increases the core genome size to

266    3,809,101 bp, consisting of 823,064 whole-genome SNPs, with 664,902 CDS SNPs. With such a

267    dramatic increase in the number of positions used for tree inference, one can observe much

268    improved discrimination among strains, with stronger bootstrap support for most ancestral nodes

269    (Supplemental Fig. S9). This analysis is a novel approach for eukaryotic genomes, where whole

270    genome SNPs were used for phylogenetic inference instead of solely relying on gene

9

271     annotations. This method allows for the robust discrimination among *S. cerevisiae* strains and

272     helps better understand the relationships among *Saccharomyces* species.

273

274     ***Whole genome PhaME analysis of a viral outbreak***

275     The *Zaire ebolavirus* outbreak that began in 2014 was rapidly characterized by large-scale

276     sequencing and assembly of *Zaire ebolavirus* genomes from several hundred patients(Gire et al.

277     2014; Carroll et al. 2015; Simon-Loriere et al. 2015). Because PhaME can provide exquisite

278     resolution among members of the same species, it could be useful in viral pathogen outbreak

279     investigations. Therefore, PhaME was applied to a total of 538 assembled genomes sequenced

280     from the 2014/2015 *Zaire ebolavirus* outbreak, as well as 20 reference genomes sequenced from

281     previous outbreaks (Supplemental Table S4). PhaME calculated 17,561 bp as the core genome

282     size (which approaches the average *Zaire ebolavirus* genome size of ~18,959 bp) with 1399 core

283     SNP positions, of which 928 are found in annotated coding regions. The PhaME *Zaire*

284     *ebolavirus* phylogeny clearly distinguishes the previous outbreaks from the 2014 outbreak

285     (Supplemental Fig. S10). Additionally, this topology is well supported and consistent with the

286     previous trees proposed from this outbreak(Gire et al. 2014; Carroll et al. 2015; Simon-Loriere et

287     al. 2015). As previously observed, the *Zaire ebolavirus* strains recovered from Guinea and Sierra

288     Leone are interspersed throughout the tree yet form some distinct lineages(Carroll et al. 2015;

289     Simon-Loriere et al. 2015).

290

291     ***Analyzing raw metagenomic reads with PhaME***

292     Outbreaks such as the case of *Zaire ebolavirus* provide a scenario where assembly of genomes is

293     often the first step in analysis due to the fact that obtaining a pure isolate may be difficult or take

294     longer than desired for rapid analysis. Since PhaME can accurately identify where raw read

295     datasets belong in a phylogenetic tree (as shown above for reads sequenced from pure cultures),

296     we tested PhaME's ability to accurately place a known infectious agent within a phylogeny,

297     using reads derived from complex metagenomic samples. We hypothesize that a dominant clonal

298     pathogen infecting a host will be accurately placed within a phylogeny due to the read mapping

299     and SNP calling strategy in PhaME. To our knowledge, this is the first demonstrated use of raw

300     shotgun metagenomic reads (or any metagenomic data) as an input to phylogenetic

301     reconstruction software. We investigated samples taken from individuals afflicted during the

302    2014 *Zaire ebolavirus* outbreak, as well as a clinical fecal sample from a US patient having

303    returned from Germany during the 2011 *stx2*-positive Enteroaggregative *E. coli* (StxEAggEC)

304    outbreak.

305

306    For *Zaire ebolavirus* metagenomes, we selected two 2014 datasets from Sierra Leone (both have

307    few to no reads that can be matched to human sequences), one (SRX674125) with 3,505,216

308    reads, of which ~62.2% map to the *Zaire ebolavirus Mayinga 1976* reference sequence

309    (2,179,715 reads; 8,861.13 average fold coverage), and another (SRX674271) with 929,604

310    reads, of which only 0.3% can be mapped to the same reference (2,850 reads; 9.45 average fold

311    coverage). In both cases, assembled genomes derived from these two datasets are available in

312    Genbank (Supplemental Table S4) and were also used in the PhaME *Zaire ebolavirus* tree

313    (Supplemental Fig. S10). Adding the two full raw metagenome datasets to the *Zaire ebolavirus*

314    tree took 1 hour 21 minutes, did not impact the core genome or SNP statistics (Table 1), and

315    were both placed within the 2014 outbreak lineage. The sample with high *Zaire ebolavirus* load

316    was positioned very closely to the sample-matched assembled genome (Figure 4, Supplemental

317    Fig. 11). The other sample was placed as the immediate outgroup to all other 2014 strains, likely

318    stemming from the paucity of data within this dataset. However, these results highlight the power

319    of PhaME to phylogenetically characterize a target organism that comprises only a minute

320    fraction of a complex sample.

321

322    In the context of metagenomic data, the ability to accurately phylogenetically place a target

323    genome requires that a dominant clonal member of the target lineage be present in the sample in

324    order to accurately identify SNPs. With *E. coli* as a commensal resident within the human gut,

325    we tested the ability of PhaME to analyze a fecal sample derived from a patient suspected to be

326    infected with the 2011 StxEAggEC strain. A large (252,926,569 reads) fecal sample dataset was

327    mapped to 52 *E. coli* strains to evaluate the distribution of reads among available genomes.

328    While all genomes recruited reads from this dataset (total of 1.74 million reads), the dominant

329    organism within the sample appeared to belong to the StxEAggEC outbreak clade (Supplemental

330    Fig. S12).  Using the existing *E. coli*+*Shigella* PhaME phylogeny, the addition of the fecal

331    sample's raw shotgun metagenomic reads reduced the *E. coli* core genome size, as well as the

332    core SNPs (Table 1). The dominant *E. coli* within the fecal sample was accurately placed in the

333    inferred tree within the StxEAggEC phylogroup B1 (Supplemental Fig. S13), further supporting

334    the use of PhaME to establish the phylogenetic placement of infectious disease outbreak strains,

335    even when in the presence of less abundant commensal strains of the same species.

336

337    **Molecular analyses and signs of positive, negative and neutral selection**

338    Identifying SNPs found in coding regions enables further molecular evolutionary analyses –a

339    post-phylogeny option that is provided in PhaME. PhaME uses the HyPhy package Branch-Site

340    random effects likelihood (REL) model for detecting episodic diversifying selection on all genes

341    containing at least 1 SNP, as well as on a concatenation of all these genes.

342

343    The *Escherichia* and *Shigella* core genome was used to illustrate the application of molecular

344    evolutionary analyses with PhaME. A total of 324 genes were found to contain at least one SNP,

345    of which only 52 genes were detected to have lineages that showed a statistically significant

346    signal of diversifying selection (Supplemental Table S5). Only four genes showed signs of

347    positive selection, three of which were only found in the avian pathogenic *E. coli* within the tree.

348    Because the size of the core genome decreases with the addition of more genomes, the core

349    genome becomes increasingly enriched in essential genes and depleted in accessory genes. The

350    small size of the core (~5% of the average genome size) and the high number of essential genes

351    within the core genome may partially explain why so few genes show signs of diversifying

352    selection. We examined the ability of PhaME to explore a single lineage within the larger

353    phylogeny, consisting of a subset clade of 20 *Escherichia* and *Shigella* genomes. In this case, a

354    total of 2513 genes were found to contain at least one SNP and, in 111 instances, a gene showed

355    a statistically significant signal of diversifying selection, either in a single strain or within a

356    monophyletic lineage comprising multiple strains. While nine genes of varied function showed

357    signs of positive selection (Supplemental Fig. S14), most genes are under negative selection,

358    including a number of genes involved in metabolism and housekeeping functions (Supplemental

359    Table S6).

360

361    **DISCUSSION**

362    With the rapidly growing number of available genomes and NGS read datasets, it is becoming

363    increasingly important to have appropriate analysis tools that can deal with both assembled

12

364    contigs (or complete genomes), as well as raw sequencing data. Several considerations in

365    developing tools include: 1) the ability to handle short reads with errors while still producing

366    accurate results, and 2) the ability to process large amounts of data in reasonable timeframes. It is

367    also becoming increasingly important that tools be designed modularly to accommodate different

368    research goals and that the tools be widely applicable. Here, we describe a new phylogenetic

369    tool, PhaME, that can rapidly process hundreds of genomes and/or raw reads to obtain highly

370    robust, whole genome SNP phylogenetic trees, and to estimate molecular evolution along

371    lineages. This tool is novel due to its ability to: 1) incorporate both raw read datasets (including

372    metagenomic) and genome assemblies, 2) uniquely combine an automatic core genome search,

373    SNP identification and phylogenetic tree generation, and 3) assess the selective pressure along

374    lineages.

375

376    We have tested PhaME using viral, eukaryotic, and bacterial genomes, and have constructed

377    trees that include many hundreds of genomes of the same genus, as well as different but related

378    genera. Using these examples from across the tree of life, we have demonstrated the ability to

379    rapidly process up to almost 1 TB of data and to construct highly robust phylogenies, some of

380    which have never before been done at this scale. When given an annotation for one of the

381    reference genomes, we have also shown the ability to perform molecular evolutionary analyses.

382    Although inclusion of this option significantly increases the runtime and memory required, it can

383    be used to assess positive selection in distinct lineages, and can lead to hypotheses based on

384    biologically relevant genomic signals.

385

386    We have shown PhaME's ability to construct phylogenies with the inclusion of raw

387    metagenomic data, a feature unique to PhaME. In two examples using sequenced clinical

388    samples from independent outbreaks, one with *Zaire ebolavirus* and one with *E. coli*, the

389    placement of the sample within the context of the other strains use in the phylogeny support

390    assumptions regarding the identity of the pathogen found within the samples. In the *Zaire*

391    *ebolavirus* example, with as little as 0.3% of the sample reads belonging to the actual pathogen,

392    the placement within the tree lies within the correct clade of the published outbreak genomes.

393    With *E. coli*, despite having a dominant presence of the StxEAggEC within the sample,

394    significant amounts of data could be best matched with other (commensal) *E. coli*, presumably

13

395    originating from the patient's normal microflora. In this case, the robust placement in the tree is

396    due the dominance of the target strain and the algorithms used to construct the alignment and

397    phylogeny. Given the growing trend in sequencing metagenomes from ill individuals, this

398    particular application of PhaME may be highly useful in a clinical setting.

399

400    We have included within PhaME the ability to incrementally add samples to a previously

401    constructed core alignment (and tree), which allows for rapid analysis of additional samples, and

402    have similarly included the ability to rapidly recompute the core genome and derived

403    phylogenies using subsets of the input genomes. While phylogenetic analysis has traditionally

404    required annotated genes, here we provide a highly automated process for today's genomic era,

405    agnostic of the input sequencing data that enables the construction of whole genome core

406    alignment, phylogenetic, and molecular evolutionary analyses within a single tool.

407

408    **METHODS**

409    **Modular Design and PhaME Overview**

410    We present a tool for Phylogenetic and Molecular Evolution Analyses (PhaME) that can take

411    raw NGS reads or assembled contigs representing draft or complete genomes, align the core

412    conserved sections among the genomes, identify all SNPs (both coding and non-coding), infer a

413    phylogenetic tree, and perform evolutionary analyses, such as identifying signals of positive

414    selection. PhaME is a Perl program that incorporates several other open source software,

415    including the MUMmer package NUCmer(Kurtz et al. 2004), Bowtie2 2.1.0(Langmead and

416    Salzberg 2012), BWA 0.7.5a(Li and Durbin 2009), SAMtools 0.1.18(Li et al. 2009),

417    HyPhy(Pond et al. 2005), PAML(Yang 2007), jModelTest(Darriba et al. 2012),

418    RAxML(Stamatakis 2006), FastTree(Price et al. 2010), R package APE(Popescu et al. 2012).

419    The overarching architecture of the PhaME analysis pipeline is outlined in Figure 1. PhaME is

420    run using a control file (SNPphy.ctl) that can be modified by the user, and requires a minimum of

421    one reference genome in FASTA format, consisting of one or more contigs (which can also be

422    complete chromosomes, megaplasmids, plasmids, etc.). Additional genomes in the following

423    formats can also be included in the analyses: finished genomes and draft contigs as FASTA files,

424    and raw next generation sequencing reads in FASTQ format. The directory structure and output

425    files created by the PhaME analysis pipeline are outlined in Supplemental Figure S15.

14

426

427     The main outputs of PhaME include all pairwise contig alignments, the core genome alignment

428     (all orthologous positions conserved among all input genomes), the final alignment of all

429     positions with one or more SNPs (a subset of the core), maximum likelihood tree(s), text files

430     summarizing the number of SNPs in pairwise comparisons, the positions of all SNPs in all input

431     genomes, and information on whether these SNPs alter a codon and its associated amino acid

432     (for a full list of output files, see Supplemental Fig. S15). The molecular evolution analyses are

433     performed on each gene that contains a SNP (when a GFF annotation file is provided for the

434     reference) and are presented in a series of files per gene as well as in a summary file for all

435     genes. Additional information on PhaME software can be found at https://github.com/LANL-

436     Bioinformatics/PhaME.

437

438     **Whole genome alignments and SNP discovery from genomes, contigs and reads**

439     Any input contigs (draft or complete genomes) are initially subjected to self-comparisons using

440     NUCmer in order to remove duplicated regions or other highly similar 'repetitive' elements

441     (length $\geq$ 100; identity 95) that could give misleading alignment results. Whole and draft genome

442     sequences then undergo pairwise whole genome alignment (in all combinations) using NUCmer.

443     Gap regions, consisting of unaligned segments $\geq$1 nucleotide are considered evolutionarily

444     uninformative and are removed. Input raw read datasets (either single or paired-end) are then

445     individually aligned to a selected reference using BowTie2 (default) or BWA. The alignment

446     results are parsed using SAMtools and perl scripts to group identified SNPs by shared genomic

447     location within the selected reference. An orthologous SNP alignment is created for each

448     genome, contig, and/or read set, and contains the nucleotides at all positions that are found in all

449     genomes, and where at least 1 genome differs at that position. Given an annotation file in GFF

450     format, the pipeline can also distinguish SNPs present within coding sequence (CDS) regions

451     from those present in intergenic spaces. The CDS SNPs are used to build a separate phylogenetic

452     tree and will also be used for downstream evolutionary analyses (when selected). The SNPs

453     identified in all pairwise genome alignments, as well as those identified using reads mapped to

454     one of the genomes are available as text files. In addition, pairwise SNP profiles for i) the core

455     genome, as well as for ii) the core coding genome and iii) the core intergenic genome (when

15

456    annotation is provided) are also available. These SNP matrices allow for rapid recalculation of

457    the core sequences for any subset of genomes, and then reconstruction of subtrees.

458    The user is provided a proximity cut-off option that excludes SNPs that are located within the

459    proximity. This feature may prove useful if horizontal transfer has mobilized multiple SNPs from

460    one genome to another, confusing the evolutionary signal(Foster et al. 2009; Croucher et al.

461    2011).  If the user does not select a reference genome for the above analyses, one will be selected

462    randomly from among those with GFF files. If no GFF files are present, then a reference is

463    chosen at random, and the optional CDS-specific modules (including the molecular evolutionary

464    analyses) will not be conducted. Example output files for each of the steps are available from

465    https://github.com/LANL-Bioinformatics/PhaME.

466

**Phylogenetic Reconstruction**

468    The SNP alignment, consisting of all SNP positions present in any genome within the identified

469    core genome, is used to construct a phylogenetic tree. If a GFF annotation file is provided, an

470    additional tree can be generated from the subset of SNPs found only within coding sequences, or

471    only within the intergenic spaces. The phylogenetic trees are inferred with using FastTree

472    (default) and/or the RAxML-HPC maximum-likelihood method. When RAxML is selected,

473    jModelTest is first run to determine the best substitution model to use when inferring the trees.

474    Both RAxML and FastTree produce a newick file that can be viewed using a display tool such as

475    FigTree (http://tree.bio.ed.ac.uk/software/figtree/), Dendroscope (http://ab.inf.uni-

476    tuebingen.de/software/dendroscope/)(Huson and Scornavacca 2012), Archeaopteryx

477    (https://sites.google.com/site/cmzmasek/home/software/archaeopteryx), EvolView(Zhang et al.

478    2012), iTOL (http://itol.embl.de)(Letunic and Bork 2011), and Phylodendron

479    (http://iubio.bio.indiana.edu/soft/molbio/java/apps/trees/).  For convenience, we have included a

480    basic tree view within a PDF file using jphylo (http://sourceforge.net/projects/jphylo/).

481

**Molecular evolutionary analyses**

483    When selected, the first step in molecular evolutionary analyses is to root the above generated

484    tree using the APE package in R. Using the reference GFF file, all homologous genes containing

485    SNPs are used to test for positive or negative selection through the implementation of methods

486    within the HyPy or PAML packages. Both packages can test for the presence of positively

16

487    selected sites and lineages by allowing the $d_N/d_S$ ratio ($\omega$) to vary among sites and lineages. The

488    branch-site REL model in the HyPhy package (hyphy.org) is used to detect instances of episodic

489    diversifying and purifying selection. If PAML is selected, the M1a-M2a and M7-M8 nested

490    models are implemented (molecularevolution.org/software/phylogenetics/paml). In this latter

491    case, the likelihood ratio test between the null models (M1a and M8) and the alternative model

492    (M2a and M7) at a significance cutoff of 5% provides information on how the genes are

493    evolving.  The results for each gene are then summarized in a table containing information on

494    whether the gene is evolving under positive, neutral, or negative selection, along with $p$-values.

495    Using HyPhy, a phylogenetic tree with evolutionary information on each branch is generated as a

496    postscript file (for each gene as well as for the concatenation of all genes). With PAML, the

497    newick tree file is modified to incorporate evolutionary information on any branch that is found

498    to be under positive selection. We opted to provide PAML as an option, however recommend

499    using HyPhy (set as default) for genome-scale projects due to its speed and concise output.

500    **Data Access**
501    GenBank accession numbers for the sequencing data and genomes used in this study can be
502    found in Supplemental Tables 1-4. The PhaME software together with documentation can be
503    found at https://github.com/LANL-Bioinformatics/PhaME.
504
511    **Author contributions**: P.S.G.C. conceived the study. S.A.A. was responsible for algorithm
512    design and bioinformatics analyses. C.L. was responsible for read mapping and alignment to
513    references. P.L. and C.L. were responsible for making the pipeline available online. K.W.D.
514    generated the flowchart for pipeline design. P.S.G.C. and S.A.A. interpreted the data and wrote
515    the manuscript with input from the other authors.
516

517

518

519

520

521

522

17

523

524   **Figure legends**

525   **Figure 1. PhaME analysis pipeline.** The PhaME analysis pipeline is able to identify SNPs

526   from complete, assembled, and read datasets and infer a phylogenetic tree. If assembled genomes

527   are provided, NUCmer is used to identify repeats and perform pairwise alignments. Bowtie2 or

528   BWA is used to map reads to one of the repeat-masked genomes. The SNP and gap coordinates

529   are used to generate whole-genome SNP alignments. If an annotation file is provided, a separate

530   alignment consisting of SNPs only found in the CDS regions is generated. RAxML or FastTree

531   phylogenies are constructed using the SNP alignments. If specified, PAML or HyPhy packages

532   can be used to test for episodic diversifying selection.

533

534   **Figure 2. Inter-genus phylogeny of *Escherichia*, *Shigella* and *Salmonella* strains.** Whole-

535   genome SNPs from 413 *Escherichia*, *Shigella*, and *Salmonella* genomes were used to construct a

536   maximum likelihood phylogenetic tree. *Shigella* and *E. coli* strains cluster together into

537   previously identified phylogroups (colored, collapsed clades), with the *Salmonella* genomes as a

538   distinct outgroup. Stars (*) represent bootstrap values under 75, all other clades are supported by

539   values >75. Entries with _c represent contig datasets, the remaining datasets are finished

540   genomes. See Supplemental Figure S1 for a fully detailed phylogenetic tree of this group. Scale

541   = 0.5 substitutions per base.

542

543   **Figure 3. Phylogenetic tree of the *Burkholderia* genus using reads, contigs and finished**

544   **genomes.** Maximum likelihood phylogeny of 217 *Burholderia* genome datasets using PhaME.

545   Collapsed branches represent various clades of the mentioned groups. Stars (*) represent

546   bootstrap values under 75, all other clades are supported by values >75. Supplemental Figure S4

547   displays the detailed phylogenetic tree with all 217 entries, and demonstrates the placement of

548   full genomes, assemblies and reads, which generally cluster together. Supplemental Figure S5

549   shows the relationships among the *Bcc* genomes together with outgroups. Supplemental Figure

550   S7 provides a more detailed view of the placement of reads with contigs and genomes within the

551   *B. pseudomallei* and *B. mallei* clade. Entries with _f, _c, _r represent: _f, finished genomes; _c,

552   contigs from assembled genomes; r, raw reads. Scale = 0.09 substitutions per base.

553

18

554     **Figure 4. Accurate placement of metagenomic samples within the *Zaire ebolavirus***

555     **phylogeny.** PhaME was used to generate a maximum likelihood phylogeny of 560 *Zaire*

556     *ebolavirus* genomes (20 reference genomes, 538 assembled genomes from 2014 *Zaire ebolavirus*

557     outbreak, and two raw metagenomic read datasets from individuals suspected of having been

558     infected in this outbreak). Collapsed branches represent various clades of *Zaire ebolavirus* that

559     are consistent with those reported previously(Gire et al. 2014)·(Carroll et al. 2015). Supplemental

560     Figure S11 displays the detailed phylogenetic tree with all 560 entries. Entries with _r represent

561     raw metagenomic read datasets. Entries labeled as: [$] represent data from reference(Gire et al.

562     2014); [¥] represent data from reference(Carroll et al. 2015); [§] represent data from

563     reference(Simon-Loriere et al. 2015). Scale = 0.02 substitutions per base.
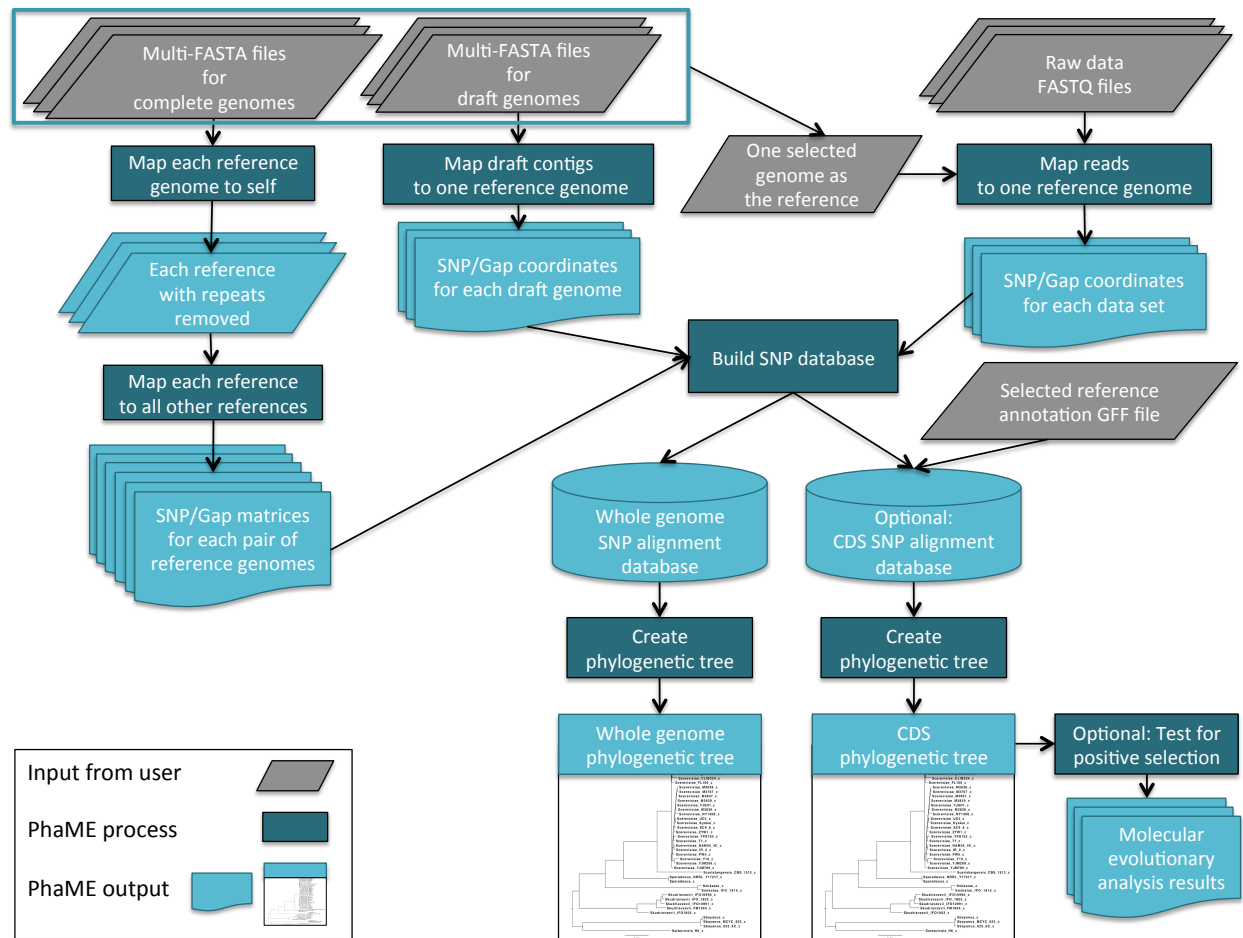
564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

19

585    Figure 1

588    Figure 2



589

590

591     Figure 3



592

593

594    Figure 4



595

596

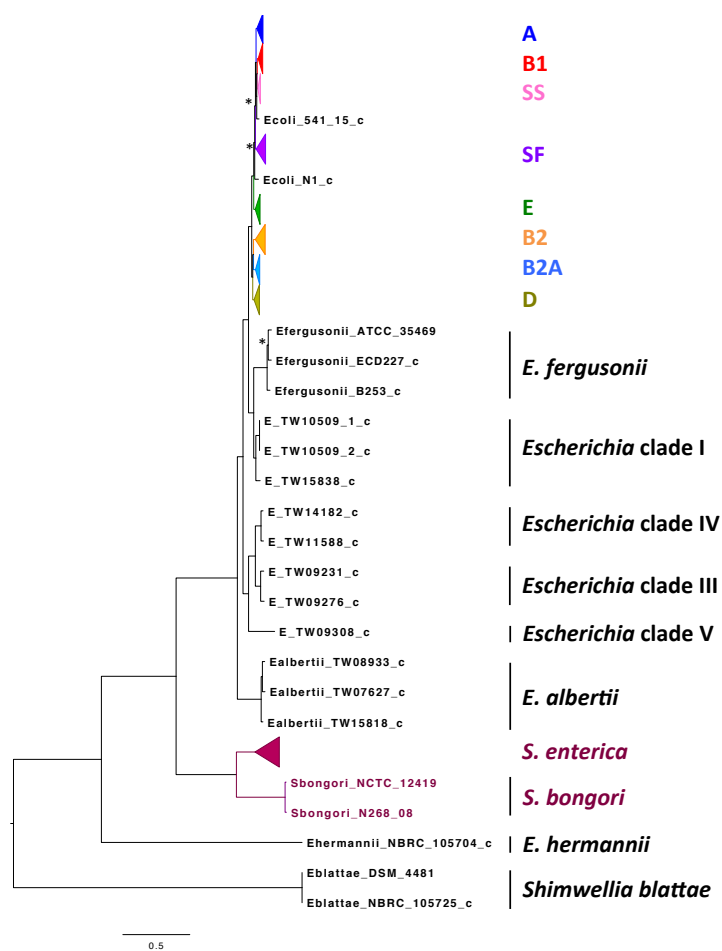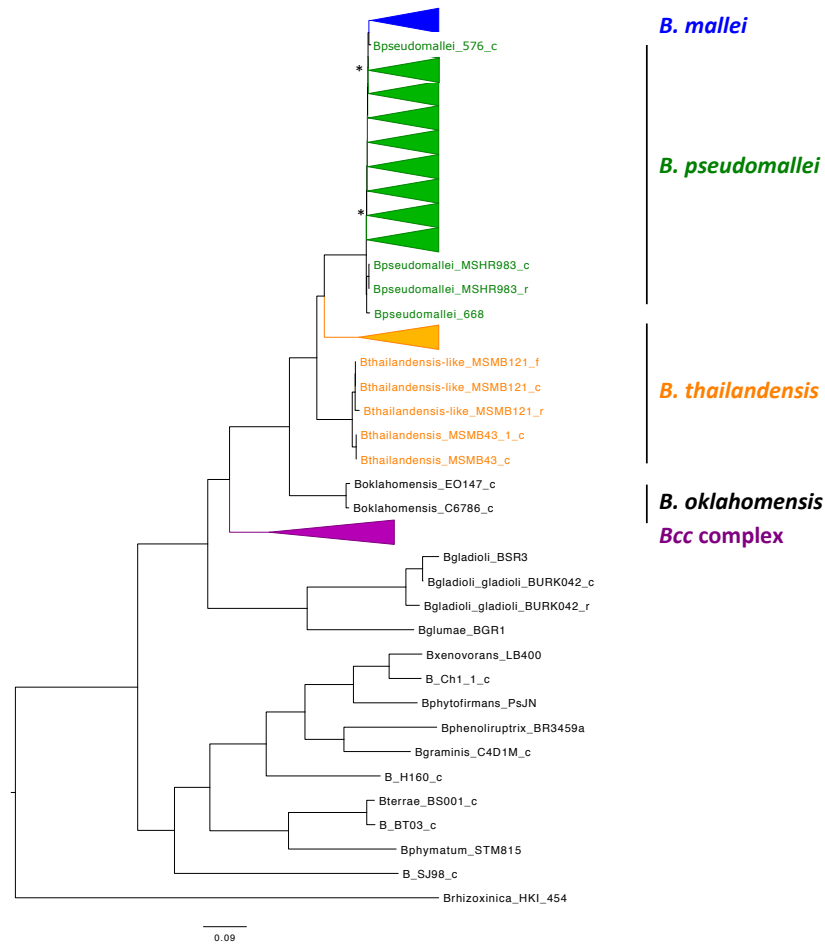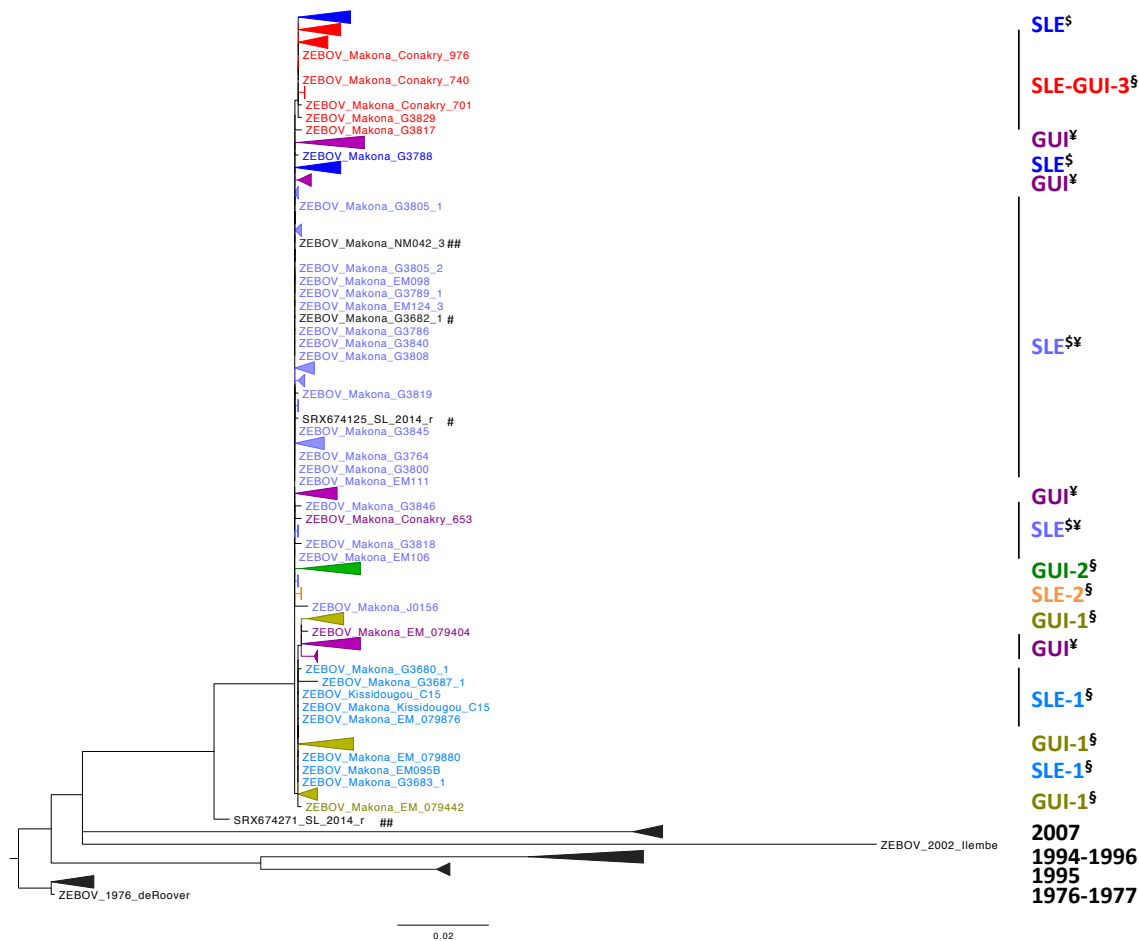## Tables

Table 1. Summary of genomes, assemblies and read datasets together with core genome statistics.

| | Escherichia/Shigella | Ecoli/Shigella | Metagenome addition | Escherichia/Shigella/Salmonella | True Escherichia/Shigella | Burkholderia | Bpseudomallei/Bmallei | Burkholderia BCC | Burkholderia BCC/pseudomallei/xenovorans | Saccharomyces | Scerevisiae | Zaire ebolavirus |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Number of NCBI genomes (Complete and Draft) | 302 | 291 | 294 | 413 | 308 | 163 | 118 | 25 | 27 | 77 | 45 | 558 |
| Number of NGS read files | 0 | 0 | 1 | 0 | 0 | 31 | 49 | 0 | 0 | 2 | 0 | 2 |
| Total genomes | 302 | 291 | 295 | 413 | 308 | 194 | 167 | 25 | 27 | 79 | 45 | 560 |
| Average genome size | 4966643 | 4999028 | 4966643 | 4783138 | 4999028 | 6957651 | 6901660 | 7367266 | 7544548 | 12071326 | 12071326 | 18613 |
| Total gap length | 4904879 | 3331000 | 3620525 | 4986576 | 3686851 | 6279803 | 3956942 | 5478178 | 6071363 | 10929991 | 8262225 | 1398 |
| Total repeats | 180894 | 180894 | 180894 | 180894 | 180894 | 226210 | 226210 | 113966 | 113966 | 874683 | 874683 | 0 |
| Core genome size | 249983 | 1823862 | 1534337 | 168191 | 1468011 | 1028251 | 3351112 | 1800938 | 1207753 | 1141335 | 3809101 | 17561 |
| Total core SNPs | 73640 | 541834 | 371060 | 49337 | 461639 | 72968 | 89399 | 391765 | 274048 | 463 | 823064 | 1399 |
| CDS SNPs | 71341 | 502760 | 347746 | 47813 | 433890 | 72642 | 71479 | 380062 | 269380 | 463 | 664902 | 928 |
| % core of the genome | 5.03 | 36.48 | 30.89 | 3.52 | 29.37 | 14.78 | 48.56 | 24.45 | 16.01 | 9.45 | 31.55 | 94.35 |
| % total core SNPs of core | 29.46 | 29.71 | 24.18 | 29.33 | 31.45 | 7.10 | 2.67 | 21.75 | 22.69 | 0.04 | 21.61 | 7.97 |
| % CDS SNPs of total SNPs | 96.88 | 92.79 | 93.72 | 96.91 | 93.99 | 99.55 | 79.96 | 97.01 | 98.30 | 100.00 | 80.78 | 66.33 |
| Wallclock time | | | 0:02:06:28 | 1:02:39:06 | | 0:14:13:44 | | | | 0:15:32:36 | 0:1:26:48 | |
| User time | | | 0:16:57:07 | 19:13:32:22 | | 1:23:09:59 | | | | 0:14:18:05 | 0:08:25:48 | |
| System time | | | 0:01:40:06 | 0:14:29:11 | | 0:00:59:09 | | | | 15:17:0:35 | 0:01:09:19 | |
| Maximum vmem | | | 9.892G | 50.720G | | 66.394G | | | | 51.405G | 8.479G | |

**REFERENCES**

Ahmed SA, Awosika J, Baldwin C, Bishop-Lilly KA, Biswas B, Broomall S, Chain PS, Chertkov O, Chokoshvili O, Coyne S et al. 2012. Genomic comparison of Escherichia coli O104:H4 isolates from 2009 and 2011 reveals plasmid, and prophage heterogeneity, including shiga toxin encoding phage stx2. *PLoS One* **7**: e48228.

Baldwin A, Mahenthiralingam E, Drevinek P, Vandamme P, Govan JR, Waine DJ, LiPuma JJ, Chiarini L, Dalmastri C, Henry DA et al. 2007. Environmental Burkholderia cepacia complex isolates in human infections. *Emerg Infect Dis* **13**: 458-461.

Ben-Ari G, Zenvirth D, Sherman A, Simchen G, Lavi U, Hillel J. 2005. Application of SNPs for assessing biodiversity and phylogeny among yeast strains. *Heredity (Edinb)* **95**: 493-501.

Brisse S, Verduin CM, Milatovic D, Fluit A, Verhoef J, Laevens S, Vandamme P, Tummler B, Verbrugh HA, van Belkum A. 2000. Distinguishing species of the Burkholderia cepacia complex and Burkholderia gladioli by automated ribotyping. *J Clin Microbiol* **38**: 1876-1884.

Carroll MW Matthews DA Hiscox JA Elmore MJ Pollakis G Rambaut A Hewson R Garcia-Dorival I Bore JA Koundouno R et al. 2015. Temporal and spatial analysis of the 2014-2015 Ebola virus outbreak in West Africa. *Nature* **524**: 97-101.

Chain PS, Denef VJ, Konstantinidis KT, Vergez LM, Agullo L, Reyes VL, Hauser L, Cordova M, Gomez L, Gonzalez M et al. 2006. Burkholderia xenovorans LB400 harbors a multi-replicon, 9.73-Mbp genome shaped for versatility. *Proc Natl Acad Sci U S A* **103**: 15280-15287.

Chaudhuri RR, Henderson IR. 2012. The evolution of the Escherichia coli phylogeny. *Infect Genet Evol* **12**: 214-226.

Coenye T, Vandamme P. 2003. Diversity and significance of Burkholderia species occupying diverse ecological niches. *Environ Microbiol* **5**: 719-729.

Coenye T, Vandamme P, Govan JR, LiPuma JJ. 2001. Taxonomy and identification of the Burkholderia cepacia complex. *J Clin Microbiol* **39**: 3427-3436.

Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J, van der Linden M, McGee L, von Gottberg A, Song JH, Ko KS et al. 2011. Rapid pneumococcal evolution in response to clinical interventions. *Science* **331**: 430-434.

Daligault HE, Davenport KW, Minogue TD, Bishop-Lilly KA, Broomall SM, Bruce DC, Chain PS, Coyne SR, Frey KG, Gibbons HS et al. 2014. Whole-genome assemblies of 56 burkholderia species. *Genome Announc* **2**.

Darriba D, Taboada GL, Doallo R, Posada D. 2012. jModelTest 2: more models, new heuristics and parallel computing. *Nat Methods* **9**: 772.

Davis S, Pettengill JB, Luo Y, Payne J, Shpuntoff A, Rand H, Strain E. 2015. CFSAN SNP Pipeline: an automated method for constructing SNP matrices from next-generation sequence data. *PeerJ Computer Science* **1**: e20.

Faison WJ, Rostovtsev A, Castro-Nallar E, Crandall KA, Chumakov K, Simonyan V, Mazumder R. 2014. Whole genome single-nucleotide variation profile-based phylogenetic tree building methods for analysis of viral, bacterial and human genomes. *Genomics* **104**: 1-7.

Filliol I, Motiwala AS, Cavatore M, Qi W, Hazbon MH, Bobadilla del Valle M, Fyfe J, Garcia-Garcia L, Rastogi N, Sola C et al. 2006. Global phylogeny of Mycobacterium tuberculosis based on single nucleotide polymorphism (SNP) analysis: insights into tuberculosis evolution, phylogenetic accuracy of other DNA fingerprinting systems, and recommendations for a minimal standard SNP set. *J Bacteriol* **188**: 759-772.

667  Foster JT, Beckstrom-Sternberg SM, Pearson T, Beckstrom-Sternberg JS, Chain PS, Roberto FF,
668      Hnath J, Brettin T, Keim P. 2009. Whole-genome-based phylogeny and divergence of the
669      genus Brucella. *J Bacteriol* **191**: 2864-2870.
670  Gardner SN, Hall BG. 2013. When whole-genome alignments just won't work: kSNP v2
671      software for alignment-free SNP discovery and phylogenetics of hundreds of microbial
672      genomes. *PLoS One* **8**: e81760.
673  Girault G, Thierry S, Cherchame E, Derzelle S. 2014. Application of High-Throughput
674      Sequencing: Discovery of Informative SNPs to Subtype *Bacillus anthracis*. *Advances in*
675      *Bioscience and Biotechnology* **05**: 669-677.
676  Gire SK, Goba A, Andersen KG, Sealfon RS, Park DJ, Kanneh L, Jalloh S, Momoh M, Fullah
677      M, Dudas G et al. 2014. Genomic surveillance elucidates Ebola virus origin and
678      transmission during the 2014 outbreak. *Science* **345**: 1369-1372.
679  Godoy D, Randle G, Simpson AJ, Aanensen DM, Pitt TL, Kinoshita R, Spratt BG. 2003.
680      Multilocus sequence typing and evolutionary relationships among the causative agents of
681      melioidosis and glanders, Burkholderia pseudomallei and Burkholderia mallei. *J Clin*
682      *Microbiol* **41**: 2068-2079.
683  Griffing SM, MacCannell DR, Schmidtke AJ, Freeman MM, Hyytia-Trees E, Gerner-Smidt P,
684      Ribot EM, Bono JL. 2015. Canonical Single Nucleotide Polymorphisms (SNPs) for
685      High-Resolution Subtyping of Shiga-Toxin Producing Escherichia coli (STEC) O157:H7.
686      *PLoS One* **10**: e0131967.
687  Hibbett DS, Binder M, Bischoff JF, Blackwell M, Cannon PF, Eriksson OE, Huhndorf S, James
688      T, Kirk PM, Lucking R et al. 2007. A higher-level phylogenetic classification of the
689      Fungi. *Mycol Res* **111**: 509-547.
690  Huson DH, Scornavacca C. 2012. Dendroscope 3: an interactive tool for rooted phylogenetic
691      trees and networks. *Syst Biol* **61**: 1061-1067.
692  James TY, Kauff F, Schoch CL, Matheny PB, Hofstetter V, Cox CJ, Celio G, Gueidan C, Fraker
693      E, Miadlikowska J et al. 2006a. Reconstructing the early evolution of Fungi using a six-
694      gene phylogeny. *Nature* **443**: 818-822.
695  James TY, Letcher PM, Longcore JE, Mozley-Standridge SE, Porter D, Powell MJ, Griffith GW,
696      Vilgalys R. 2006b. A molecular phylogeny of the flagellated fungi (Chytridiomycota)
697      and description of a new phylum (Blastocladiomycota). *Mycologia* **98**: 860-871.
698  Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL. 2004.
699      Versatile and open software for comparing large genomes. *Genome Biol* **5**: R12.
700  Kurtzman CP, Robnett CJ. 2003. Phylogenetic relationships among yeasts of the 'Saccharomyces
701      complex' determined from multigene sequence analyses. *FEMS Yeast Res* **3**: 417-432.
702  Lackner G, Moebius N, Partida-Martinez L, Hertweck C. 2011. Complete genome sequence of
703      Burkholderia rhizoxinica, an Endosymbiont of Rhizopus microsporus. *J Bacteriol* **193**:
704      783-784.
705  Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:
706      357-359.
707  Letunic I, Bork P. 2011. Interactive Tree Of Life v2: online annotation and display of
708      phylogenetic trees made easy. *Nucleic Acids Res* **39**: W475-478.
709  Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
710      *Bioinformatics* **25**: 1754-1760.

711   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
712         Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and
713         SAMtools. *Bioinformatics* **25**: 2078-2079.
714   Luo C, Walk ST, Gordon DM, Feldgarden M, Tiedje JM, Konstantinidis KT. 2011. Genome
715         sequencing of environmental Escherichia coli expands understanding of the ecology and
716         speciation of the model bacterial species. *Proc Natl Acad Sci U S A* **108**: 7200-7205.
717   Pamilo P, Nei M. 1988. Relationships between Gene Trees and Species Trees. *Mol Biol Evol* **5**:
718         568-583.
719   Pond SL, Frost SD, Muse SV. 2005. HyPhy: hypothesis testing using phylogenies.
720         *Bioinformatics* **21**: 676-679.
721   Popescu AA, Huber KT, Paradis E. 2012. ape 3.0: New tools for distance-based phylogenetics
722         and evolutionary analysis in R. *Bioinformatics* **28**: 1536-1537.
723   Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for
724         large alignments. *PLoS One* **5**: e9490.
725   Priest FG, Barker M. 2010. Gram-negative bacteria associated with brewery yeasts:
726         reclassification of Obesumbacterium proteus biogroup 2 as Shimwellia pseudoproteus
727         gen. nov., sp. nov., and transfer of Escherichia blattae to Shimwellia blattae comb. nov.
728         *Int J Syst Evol Microbiol* **60**: 828-833.
729   Reik R, Spilker T, Lipuma JJ. 2005. Distribution of Burkholderia cepacia complex species
730         among isolates recovered from persons with or without cystic fibrosis. *J Clin Microbiol*
731         **43**: 2926-2928.
732   Sahl JW, Matalka MN, Rasko DA. 2012. Phylomark, a tool to identify conserved phylogenetic
733         markers from whole-genome alignments. *Appl Environ Microbiol* **78**: 4884-4892.
734   Sahl JW, Schupp JM, Rasko DA, Colman RE, Foster JT, Keim P. 2015. Phylogenetically typing
735         bacterial strains from partial SNP genotypes observed from direct sequencing of clinical
736         specimen metagenomic data. *Genome Med* **7**: 52.
737   Schork NJ, Fallin D, Lanchbury JS. 2000. Single nucleotide polymorphisms and the future of
738         genetic epidemiology. *Clin Genet* **58**: 250-264.
739   Simon-Loriere E, Faye O, Faye O, Koivogui L, Magassouba N, Keita S, Thiberge JM, Diancourt
740         L, Bouchier C, Vandenbogaert M et al. 2015. Distinct lineages of Ebola virus in Guinea
741         during the 2014 West African epidemic. *Nature* **524**: 102-104.
742   Song J, Xu Y, White S, Miller KW, Wolinsky M. 2005. SNPsFinder--a web-based application
743         for genome-wide discovery of single nucleotide polymorphisms in microbial genomes.
744         *Bioinformatics* **21**: 2083-2084.
745   Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
746         thousands of taxa and mixed models. *Bioinformatics* **22**: 2688-2690.
747   Storms V, Van Den Vreken N, Coenye T, Mahenthiralingam E, LiPuma JJ, Gillis M, Vandamme
748         P. 2004. Polyphasic characterisation of Burkholderia cepacia-like isolates leading to the
749         emended description of Burkholderia pyrrocinia. *Syst Appl Microbiol* **27**: 517-526.
750   Vanlaere E, Baldwin A, Gevers D, Henry D, De Brandt E, LiPuma JJ, Mahenthiralingam E,
751         Speert DP, Dowson C, Vandamme P. 2009. Taxon K, a complex within the Burkholderia
752         cepacia complex, comprises at least two novel species, Burkholderia contaminans sp.
753         nov. and Burkholderia lata sp. nov. *Int J Syst Evol Microbiol* **59**: 102-111.
754   Walk ST, Alm EW, Gordon DM, Ram JL, Toranzos GA, Tiedje JM, Whittam TS. 2009. Cryptic
755         lineages of the genus Escherichia. *Appl Environ Microbiol* **75**: 6534-6544.

756   Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**: 1586-
757          1591.
758   Zhang H, Gao S, Lercher MJ, Hu S, Chen WH. 2012. EvolView, an online tool for visualizing,
759          annotating and managing phylogenetic trees. *Nucleic Acids Res* **40**: W569-572.
760